



HAL
open science

Exploitation de Big Data Historiques pour les Humanités Numériques : application aux données financières HBDEX

Camille Guerry, Bertrand B. Coüasnon, Aurélie Lemaitre, Sébastien Adam,
Thierry Paquet, Andres Rojas Camacho

► To cite this version:

Camille Guerry, Bertrand B. Coüasnon, Aurélie Lemaitre, Sébastien Adam, Thierry Paquet, et al..
Exploitation de Big Data Historiques pour les Humanités Numériques : application aux données fi-
nancières HBDEX. NUMÉRIQUE ET PATRIMOINE Enjeux et questions actuels, Mar 2021, Paris /
Virtual, France. hal-03828391

HAL Id: hal-03828391

<https://hal.science/hal-03828391v1>

Submitted on 25 Oct 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

NUMÉRIQUE ET PATRIMOINE

Enjeux et questions actuels

DIGITAL TECHNOLOGY AND HERITAGE - CHALLENGES AND ISSUES

11/12
MARS
2021

Conférence en ligne
Digital conference

Exploitation de Big Data Historiques pour les Humanités Numériques : application aux données financières

HBDEX

Coordinateur(s): PIERRE CYRILLE HAUTCOEUR

Partenaires : PSE, LITIS, INSA-INRIA, CAMS

Présentation : Camille GUERRY, Andres ROJAS CAMACHO

Résumé :

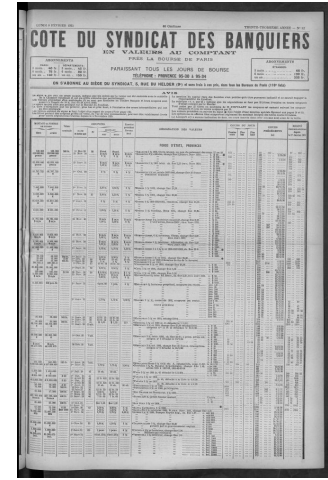
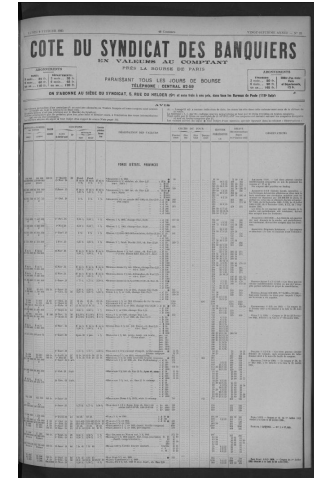
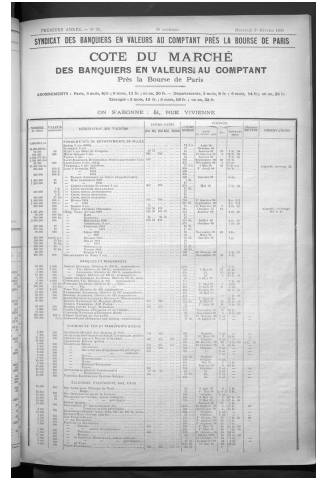
L'objectif est d'étudier le comportement à long terme des marchés financiers, en extrayant une grande quantité de données à partir d'images de documents. Nous proposons une stratégie de reconnaissance qui exploite la séquentialité des documents.



● Analyse de documents d'histoire financière ●

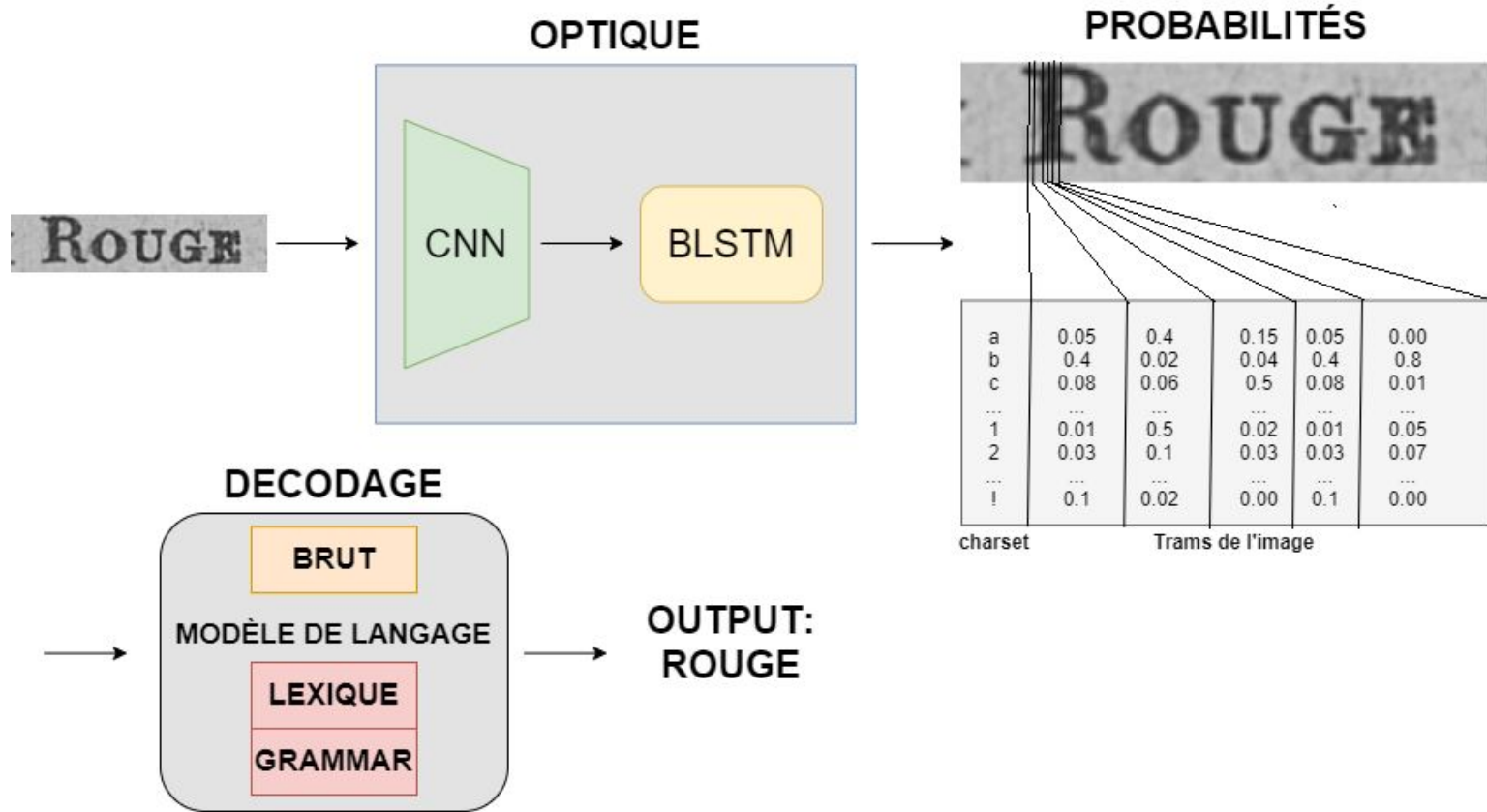
- Cotes boursières
 - **Marché de La Coulisse à Paris**
 - **1899 à 1939**
- Intérêt de l'analyse des données
 - **Modélisation des marchés financiers**
- Une des problématiques du projet
 - **Extraction des cotations dans des images de documents**
- Difficultés classiques liées aux documents anciens
 - **Documents très denses, abîmés, effacés**
 - **Impression de mauvaise qualité**

=> OCR commerciaux non satisfaisants



NOMBRE de titres	VALEUR nominale	DÉSIGNATION DES VALEURS	COURS COTÉS plus bas, plus haut, dernier	Derniers COURS	COUPONS			Dernier REVENU	OBSERVATIONS
					DATE du dernier payé	Nos	MONTANT BRUT		
MÉTALLURGIE, FABRIQUES DE MACHINES									
35 000	100	ATELIERS FRANCO-RUSSES	98	Janvier 98	3 et 4	5 fr.		
5 000	500	— DE NORD DE LA FRANCE	205	1 ^{er} Decemb. 98	15	32 fr. 60		
30 000	100	BOURKS (Usines)	150	157	Juillet 98	7	6 fr. 75		
80,875	100 r.	BRANSK (Usines de), unités	1307 50	1320	Juin 98	1807	32 rmb.		
		coupages de 10 litres	1360	Origine	1 at.			
12 000	1 000	CRANTERS DE NICOLAÏEFF	70	4 Janvier 98	16	32 fr. 70		
5 000	500	CHILERS (Hauts-fourneaux de la)	415	1 ^{er} Sept. 98	3	5 fr.		
6 000	100 d.	CLEVELAND MACHINE SCREW	310	415	Janvier 99	1	12 fr. 50		
8 000	500	CONSTR. MECAN. DU MIDI DE LA RUSSIE, actions	800	805	2			
8 000	500	— — — — — nouv.	775	1 at.			
3 000	—	— — — — — 4/10 p. de Fond.	680	Origine	1 at.			
20 000	250	CONSTRUCTIONS MÉCANIQUES (Société Française des)	263	1 ^{er} Decemb. 98	27	12 fr. 50		
20 000	500	DENAIN-ANZIN (Hauts-fourneaux, Forges et Acieries de)	1030	19 Decemb. 98	4	200 r.		
20 000	250 r.	DNIÉPROVENSNE (Société Métallurgique)	6300	6325	15 Mars 94	2 et 4			
24 000	125 r.	DONETZ (Forges et Acieries)	1025	1029	31 Aout 98	13	30 fr.		
18 000	500	DUBELANG (Hauts-fourneaux, Forges de)	1115	31 Decemb. 98	33	280 fr.		
5 000	1 000	EICH (Forges d')	5720	31 Decemb. 98	5	12 fr.		
3 200	500	ÉTABLISSEMENTS BARBIER ET C ^o	640	655	15 Février 98	2	10 fr.		

OCR : Schéma fonctionnement



OCR: Exemple de grammaire

MONTANT ou NOMBRE DE TITRES		VALEUR nominale	COUPONS			Exercice précédent Revenu brut	DESIGNATION DES VALEURS	COURS DU JOUR				CLOTURE PRÉCÉDENTE (A)	RELEVÉ des cours extrêmes depuis le 1 ^{er} Janv. 1924		
ÉMIS	ADMIS		DATE d'échéance payé	N°	MONTANT NET BRUT			Premier cours	Plus bas	Plus haut	Dernier cours				
10.000	10.000	100 fr.	5 Mars 24	70	70 fr. 66	80 fr.	—	—	—	2500	2475	2500	2475	2100	3025
45.000	45.000	100 fr.	4 Juil. 23	2	6 fr.	6 fr.	HOTELS RÉUNIS.			390		389	386	359	545
12.000	12.000	100 fr.	1 ^{er} Mars 24	1	⊕	IMMOBILIÈRE ET HOTELIÈRE DE NORMANDIE	⊕		267		271	266	260	262
20.000	8.000	500 fr.	15 Janv. 23	17	34 fr. 35	40 fr.	LA MORUE FRANÇAISE ET SÈCHERIES DE PECAMP						725	16/4/24	725
50.000	50.000	100 fr.	17 Mars 24	52	8 fr. 14	10 fr.	L'ÉPARGNE. Alimentation. Toulouse.						298	7/1/24	298
50.000	50.000	100 fr.	5 Fév. 24	9 att.	6 fr. 92	8 fr. 20	MAISON REY	⊕		151			150		180
26.000	13.000	100 fr.	15 Déc. 20	24	6 fr. 08	7 fr.	MARGARINERIE DE BÉTHUNE						74	50 26/1/24	74 50
60.000	60.000	100 fr.	1 ^{er} Déc. 23	53	3 fr. 65	4 fr. 50	OLIBET (Société des Biscuits)			245		249 50	242	240	309
90.000	90.000	100 fr.	3 Déc. 23	10	5 fr.	6 fr.	PRUDHON (J.) ET C ^{ie} (Etablissements B. R. H. R.)			353	350		353	260	427
80.000	80.000	250 fr.	12 Déc. 23	20	\$ 1	\$ 1	RAISIN DE CORINTHE	⊕					260	11/4/24	260
80.000	80.000	7 Fév. 22	6	13 dr. 25	13 dr. 25	—						43	11/4/24	42 50
25.000	25.000	100 fr.	6 Nov. 23	17	15 fr. 25	17 fr. 50	RASPAIL (Etablissements)			400			404	364	449
4.000	4.000	500 fr.	10 Juill. 23	10	50 fr.	50 fr.	RESTAURANT HENRY						1700	26/1/24	1700

Lexique répétitif

grammaire_cours_du_jours

grammaire_cloture_précédente

grammaire_relevé_cours_ext

249 50

montant, centimes

74 50 26/1/24

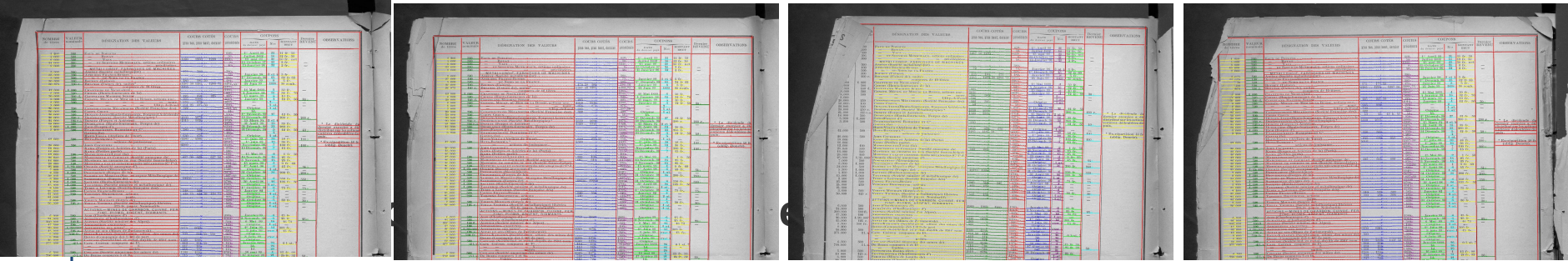
montant, centimes, date

91 25 115 ..

montant, centimes, montant, centimes

Stratégie : importance de la collection

- Exploiter les stabilités pour pallier les erreurs de reconnaissances
- Détecter automatiquement les séquences stables / points de rupture



MARCINELLE ET COULLET (Société anonyme de).....	340
MOTEURS A GAZ ET CONSTRUCTIONS MÉCANIQUES (C ^{ie} f ^{ie} d.)	160
OUGRÉ-MARIHAY (Société anonyme d').....	1170
PROVIDENCE (Forges de la).....	
SAMBRE ET MOSELLE (Soc. anonyme Métallurgique de)	
SARREBRUCK (Forges de).....	
TAGANROG (Société minière et métallurgique de).....	
TUBES A LOUVROIL (Société française des) estampillées	
USINES FRANCO-RUSSES.....	553
VERCHNY DNEPROVSK.....	
— parts.....	
VILLELONGUE (Usines Electro-Métallurgiques de).....	48
VIREUX-MOLHAIN (forges de).....	
VOLGA VICHERA (Société métallurgique) libérées	
WATTELAR-FRANCO (Usines).....	

MARCINELLE ET COULLET (Société anonyme de).....	302 50
MOTEURS A GAZ ET CONSTRUCTIONS MÉCANIQUES (C ^{ie} f ^{ie} d.)	
OUGRÉ-MARIHAY (Société anonyme d').....	
PROVIDENCE (Forges de la).....	2900
SAMBRE ET MOSELLE (Soc. anonyme Métallurgique de)	
SARREBRUCK (Forges de).....	8850
TAGANROG (Société minière et métallurgique de).....	650
TUBES A LOUVROIL (Société française des) estampillées	
USINES FRANCO-RUSSES.....	560
VERCHNY DNEPROVSK.....	27
— parts.....	
VILLELONGUE (Usines Electro-Métallurgiques de).....	
VIREUX-MOLHAIN (forges de).....	
VOLGA VICHERA (Société métallurgique) libérées	
WATTELAR-FRANCO (Usines).....	

MARCINELLE ET COULLET (Société anonyme de).....	320
MOTEURS A GAZ ET DES CONST. MÉCANIQUES (C ^{ie} f ^{ie} des)(3)	115
OUGRÉ-MARIHAY (Société anonyme d').....	
PROVIDENCE (Forges de la).....	2300
SAMBRE ET MOSELLE (Soc. anonyme Métallurgique de)	
SARREBRUCK (Forges de).....	6935
TAGANROG (Société métallurgique de).....	
... TUBES A LOUVROIL (Soc. f ^{ie} p ^{ie} la fabricat ^{ie} des) estamp.	1085
USINES FRANCO-RUSSES (Anc. Etabl. Baird).....	425
VERCHNY-DNEPROVSK (Compagnie Métallurgique de).	119
— parts.....	
VIREUX-MOLHAIN (forges de).....	
VOLGA-VICHERA (Société minière et métallurgique)	
WATTELAR-FRANCO (Usines).....	

Conclusion

- **Résultats**

- **OCR**

- Evaluation sur une page de 1899
 - > Meilleur OCR du commerce: 84.72%
 - > OCR (réseau de neurones) + Modèle de langage: 99.12%

- **Identification des titres boursiers**

- Evaluation sur 4 jours (16 images) de 1899, 1696 titres boursiers
 - > Apport de la collection : F-mesure : passe de 91,4% à 98.8 % (100% visé)

- **Perspectives**

- **Analyse des données par le CAMS**
 - **Etudes de modèles économiques par PSE**

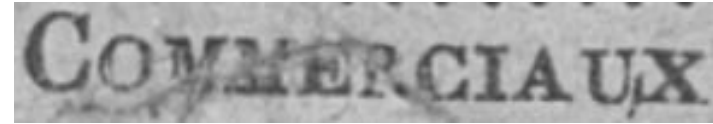
HBDEX

Slides supplémentaires

OCR: Création d'une grammaire et lexique

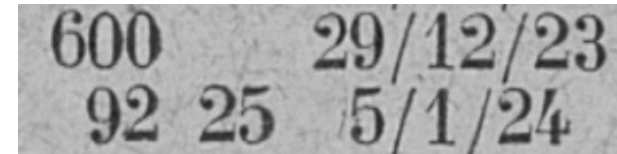
La prédiction brute se limite au visuel, elle est sensible aux imperfections dans l'image. L'utilisation d'un modèle de langage permet de corriger la sortie à l'aide d'un lexique ou d'une grammaire pré-établis.

COMmLACIAUX		COMMERCIAUX		COMMERCIAUX
		...		
		COMMUNAUTÉ		

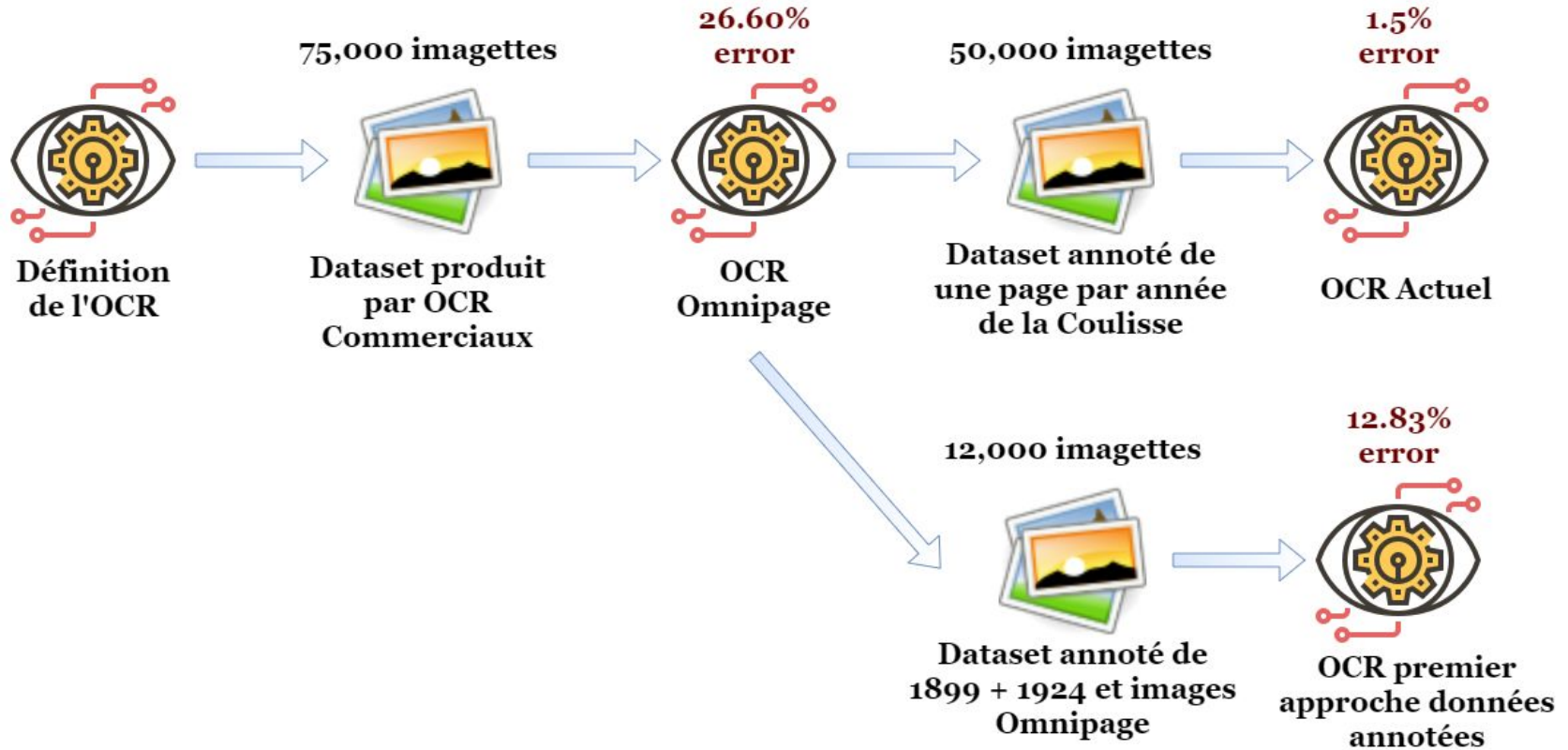


Le lexique s'utilise dans les situations où l'objet en question appartient à une liste de candidats. Pour les données qui changent un d'un jour à l'autre, il faut utiliser une grammaire pour définir un format de sortie pour avoir une prédiction correcte.

{##} {##} ##/##/## = 222 25 12/1/23



OCR : Apprentissage



OCR: Outil d'annotation

HBDEX - Côte Syndicat 1924_page8_chaque15

(Threshold : 0.9)

Exit without saving

ALIMENTATION, BRASSERIES, HOTELS,

SUCRERIES

ALIMENTATION (Société d'), en liquidation

ALIMENTATION DE PROVENCE

AU PLANTEUR DE CAIFFA, privil. 8% cum

-- ordinaires §

-- Parts

BRASSERIE DE LA COMÈTE

- DE LA MEUSE

- DE SOCHAUX

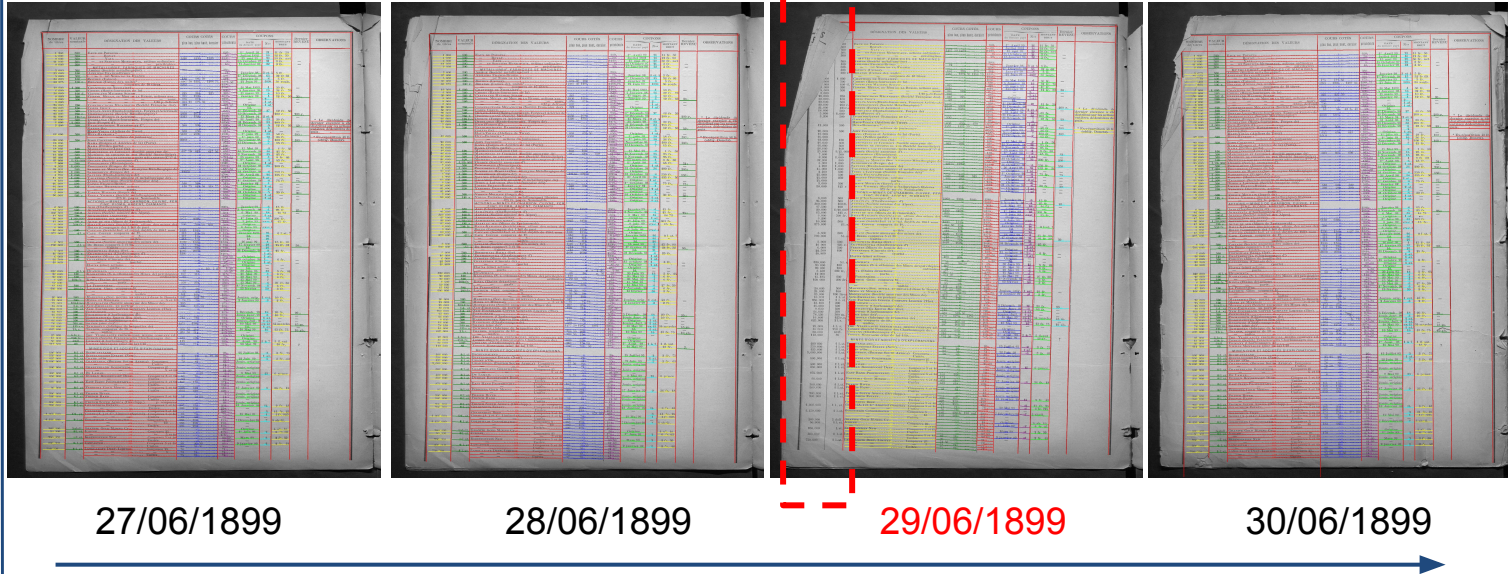
- ET TAVERNES ZIMMER §

--- Jouissance

		ALIMENTATION, BRASSERIES, HOTELS, SUCRERIES	
fr.	ALIMENTATION (Société d'), en liquidation.....
fr.	7 fr. ..	ALIMENTATION DE PROVENCE.....	245 ..
fr.	20 fr.	AU PLANTEUR DE CAIFFA, privil. 8 % cum.....	250 ..
fr.	15 fr. 40	— — ordinaires.....	502 ..
fr.	6 fr.	— — Parts.....	144 .. 140
fr.	75 fr.	BRASSERIE DE LA COMÈTE.....	1775 ..
fr.	55 fr.	— DE LA MEUSE.....	890 .. 885
fr.	60 fr.	— DE SOCHAUX.....	1059 ..
fr. 50	11 fr. 50	— ET TAVERNES ZIMMER.....	460 ..
fr. 50	6 fr. 50	— — Jouissance.....	370 ..
fr.	BRASSERIE UNIVERSELLE Jouissance.....	530 ..
fr.	25 fr.	CAFÉ-RESTAURANT AMÉRICAIN.....	880 ..
fr.	21 fr. ..	DOCKS RÉMOIS (Comptoir Gén. Alim. et App.) Unit..
fr.	7 fr.	ECO (Société Technique Alimen.), série A.....	107 50
fr.	7 fr.	— — série B.....	147 ..
lire	25 lire	*ERIDANIA (Société Industrielle)..... Unités
.....	t.c.p. 25 ..	c. 5.....

Détection de ruptures

1. Analyse d'image (grammaire et OCR) : extraction des colonnes sans a priori sur tout le corpus
2. Fiabilisation transversale : détection de ruptures et de plages stables
3. Analyse d'image : extraction des colonnes fiabilisées



Alignement de séquences

1. Analyse d'image (grammaire et OCR) : extraction des titres sans a priori sur tout le corpus
2. Fiabilisation transversale : alignement de séquences

S1	PETIT JOURNAL, 5 %	PROKHOROW	RAFFINERIE	RYKOWSKI'	Ste HOUIL. ET METALL.	TUILERIES	URIKANY, 5 %	URIKANY, 4 %
S2	PETIT JOURNAL	PROKHOROW	RAFFINERIE	RYKOWSKI	Ste HOUIL. ET METALL.	TUILERIES	URIKANY, 5 %	
S3	PETIT JOURNAL, 5 %	KHOROW	RAFFINERIE	RYKOWSKI	Ste HOUIL. ET METALL.	ATUILERIES	URIKANY, 5 %	
S4	PETIT JOURNAL, 5 %	PROKHOROW	RAFFINERIE	RIO TINTO 4 %	RYKOWSKI	Ste HOUIL.METALL.	TUILERIES	URIKANY, 5 %



S1'	PETIT JOURNAL, 5 %	PROKHOROW	RAFFINERIE	-	RYKOWSKI'	Ste HOUIL. ET METALL.	TUILERIES	URIKANY, 5 %	URIKANY, 4 %
S2'	PETIT JOURNAL	PROKHOROW	RAFFINERIE	-	RYKOWSKI	Ste HOUIL. ET METALL.	TUILERIES	URIKANY, 5 %	-
S3'	PETIT JOURNAL, 5 %	KHOROW	RAFFINERIE	-	RYKOWSKI	Ste HOUIL. ET METALL.	ATUILERIES	URIKANY, 5 %	-
S4'	PETIT JOURNAL, 5 %	PROKHOROW	RAFFINERIE	RIO TINTO 4 %	RYKOWSKI	Ste HOUIL.METALL.	TUILERIES	URIKANY, 5 %	-
	ID 1	ID 2	ID 3	ID 4	ID 5	ID 6	ID 7	ID 8	ID 9

3. Rapprochement avec le vocabulaire pour l'attribution d'un identifiant DFIH
4. Sollicitation utilisateur

1. Analyse d'image (grammaire et OCR) : extraction des titres sans a priori sur tout le corpus

9.000	500 fr.	PETIT JOURNAL, 5 o/o	507	508
10.000	500 fr.	PROKHOROW (Charbonnages de) $\frac{1}{2}$ 1/2 o/o	484	485
60.000	500 fr.	RAFFINERIE SAY, $\frac{1}{2}$ o/o	489	490
600.000 l.st.	20 l. st.	RIO TINTO $\frac{1}{2}$ o/o	467	50
3.000	500 fr.	RYKOWSKI (Charbonnages de) $\frac{1}{2}$ 1/2 o/o	467	50
5.000	500 fr.	SOCIÉTÉ DE L'HOTEL DU PALAIS (Biarritz) $\frac{1}{2}$ o/o	473	
12.000	250 fr.	Sté HOUIL. ET MÉTALL. DANS LE DONETZ MAKEWSKA, $\frac{1}{2}$ o/o		
	100 fl.	TUILERIES DU BERRY		
		URIKANY, 5 o/o		

2. Fiabilisation transversale : alignement de séquences

S1	PETIT JOURNAL, 5 %	PROKHOROW	RAFFINERIE	RYKOWSKI	Ste HOUIL. ET METALL.	TUILERIES	URIKANY, 5 %	URIKANY, 4 %
S2	PETIT JOURNAL	PROKHOROW	RAFFINERIE	RYKOWSKI	Ste HOUIL. ET METALL.	TUILERIES	URIKANY, 5 %	
S3	PETIT JOURNAL, 5 %	KHOROW	RAFFINERIE	RYKOWSKI	Ste HOUIL. ET METALL.	ATUILERIES	URIKANY, 5 %	
S4	PETIT JOURNAL, 5 %	PROKHOROW	RAFFINERIE	RIO TINTO 4 %	RYKOWSKI	Ste HOUIL.METALL.	TUILERIES	URIKANY, 5 %

S1'	PETIT JOURNAL, 5 %	PROKHOROW	RAFFINERIE	-	RYKOWSKI	Ste HOUIL. ET METALL.	TUILERIES	URIKANY, 5 %	URIKANY, 4 %
S2'	PETIT JOURNAL	PROKHOROW	RAFFINERIE	-	RYKOWSKI	Ste HOUIL. ET METALL.	TUILERIES	URIKANY, 5 %	-
S3'	PETIT JOURNAL, 5 %	KHOROW	RAFFINERIE	-	RYKOWSKI	Ste HOUIL. ET METALL.	ATUILERIES	URIKANY, 5 %	-
S4'	PETIT JOURNAL, 5 %	PROKHOROW	RAFFINERIE	RIO TINTO 4 %	RYKOWSKI	Ste HOUIL.METALL.	TUILERIES	URIKANY, 5 %	-
	ID 1	ID 2	ID 3	ID 4	ID 5	ID 6	ID 7	ID 8	ID 9

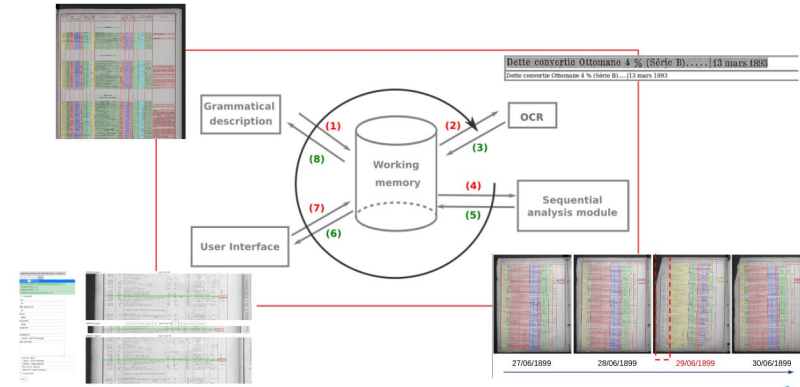
3. Rapprochement avec le vocabulaire pour l'attribution d'un identifiant DFIIH

4. Sollicitation utilisateur

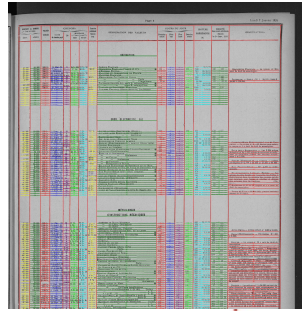
Application à la collection de listes de prix

Stratégie complète de reconnaissance

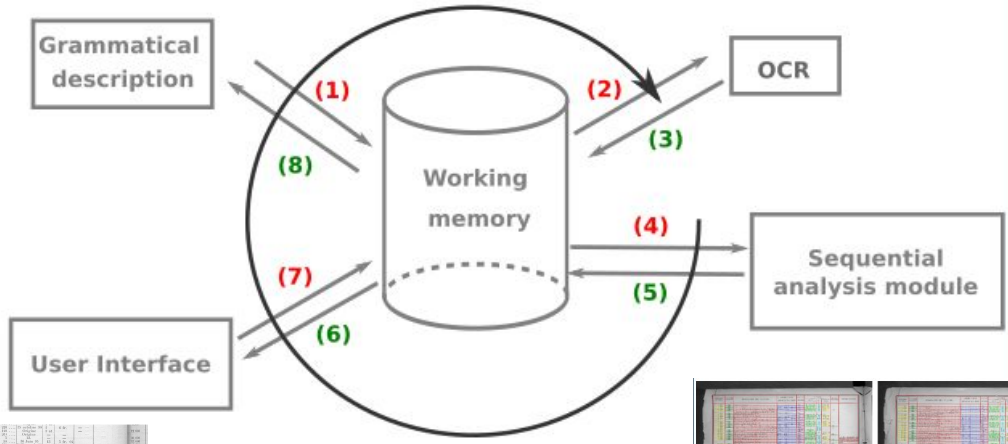
- **Itération 1 : colonnes**
 - Extraction des colonnes sans a priori sur tout le corpus
 - Fiabilisation transversale
 - Extraction des colonnes fiabilisées
- **Itération 2 : sections**
 - Extraction des sections sans a priori sur tout le corpus
 - Fiabilisation transversale
 - Extraction des sections fiabilisées
- **Itération 3 : titres**
 - Extraction des titres sans a priori sur tout le corpus
 - Fiabilisation transversale
 - Attribution d'un identifiant DFIH
 - Sollicitation utilisateur
- **Itération 4 : autres champs**
 - Fiabilisation transversale
 - Identification des devises
 - Sollicitation utilisateur
- **Production du XML**



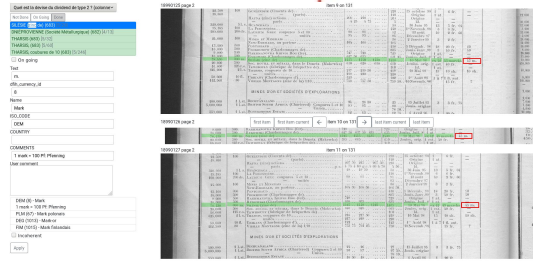
Systeme d'analyse de collection



Dette convertie Ottomane 4 % (Série B)..... | 13 mars 1893
 Dette convertie Ottomane 4 % (Série B)..... | 13 mars 1893



Visio - 27/01/2021



27/06/1899 28/06/1899 29/06/1899 30/06/1899

OCR: Résultats OCR + Modèle de langage

Test page 1899

OCR	Character Error Rate
Tesseract	31.62%
Omnipage	15.28%
OCR Omni + 1924 + 1899	11.03%
OCR Omni + 1924 + 1899 + Modèle de langage	3.60%
OCR une page par année	1.41%
OCR une page par année + Modèle de langage	0.88%

Résultats quantitatifs

Nombre d'images	Nb d'erreurs avant fiabilisation	Nb d'erreur après fiabilisation
331	14	0

Identification des titres boursiers

	Précision	Rappel	F-mesure
Pas de prise en compte de la collection (pages isolées)	0.913	0.914	0.914
Prise en compte de la collection (pas d'interaction)	0.931	0.932	0.931
Prise en compte de la collection et interactions	0.989	0.988	0.988

Evaluation sur 4 jours (16 images) de 1899

-> 1696 titres boursiers

Utilisation de OCR avec erreur de 0.88%

0% d'erreur visé suite au dernières questions

Résultats quantitatifs

Identification des titres boursiers

Evaluation sur 4 jours (16 images) de 1899
-> 1696 titres boursiers

Vocabulaire utilisé :

ensemble des titres sans les
variations dans les intitulés

Utilisation de OCR avec erreur de 0.88%

F-mesure sur titres boursiers
Sans collection : 0,914
Avec collection : 0,988

0% d'erreur visé suite
au dernières questions

	Pas de prise en compte de la collection (pages isolées)	Prise en compte de la collection Pas d'interaction	Prise en compte de la collection Après interactions
Nb de vrai positif	1550	1580	1675
Nb faux positif	147	117	13
Nb faux négatif	146	116	21
Précision	0.913	0.931	0.989
Rappel	0.914	0.932	0.988
F-mesure	0.914	0.931	0.988

- **Conception d'une stratégie d'analyse**
 - **Réalisation d'itérations successives**
 - **Construction d'un OCR adapté**
 - **Analyse transverse et fiabilisation**
 - **Applications à des éléments de granularité variable**

- **Résultats :**
 - **Moins de 1% pour l'OCR**
 - **Correction automatique des colonnes**
 - **Correction automatique des sections**
 - **Identifications des titres par fiabilisation avec la collection**
 - **F-mesure : passe de 91,4% à 98.8 % (100% visé)**