

Extraction automatique d'informations dans les annuaires commerciaux parisiens

Nathalie Abadie ⁽¹⁾, Edwin Carlinet ⁽²⁾, Joseph Chazalon ⁽²⁾, Bertrand Duménieu ⁽³⁾.

Séminaire SoDUCo - BnF, 10 novembre 2022



(1)



(2)



(3)



SODUCO ANR-18-CE38-0013

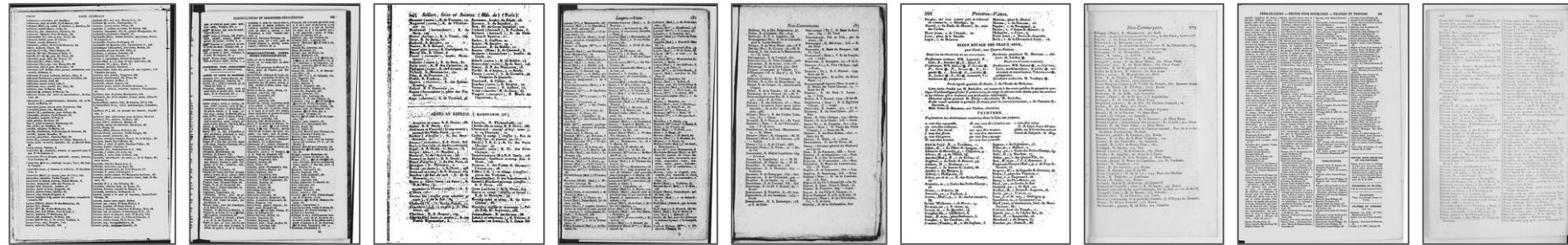


Introduction

Un corpus de **141** annuaires numérisés :

- Plusieurs centaines de milliers de pages,
- Plusieurs millions d'**entrées**, comportant le nom des commerces, leur type d'activité et leur localisation.

⇒ Une immense **source de données** pour le suivi individuel des commerces au cours du 19^e siècle...
... qu'il faut **structurer** pour permettre d'automatiser les analyses et les traitements !



Des documents numérisés aux données structurées



Lemonnyer, plomberie, r. de Bondy, 86, et
r. Bouchardon, 1.

Lemonnyer, plomberie, r. de Bondy, 86, et r. Bouchardon, 1.

PERSONNE	ACTIVITÉ	RUE	NUMÉRO	RUE	NUMÉRO
Lemonnyer,	plomberie,	r. de Bondy,	86,	et r. Bouchardon,	1.



Pages d'annuaires numérisées

Identification des entrées

Entrée d'annuaire

Transcription du texte

Texte numérique de l'entrée

Annotation du texte

Entités nommées

Structuration des entrées

Base de données intégrée

Construire un outil pour les historiens

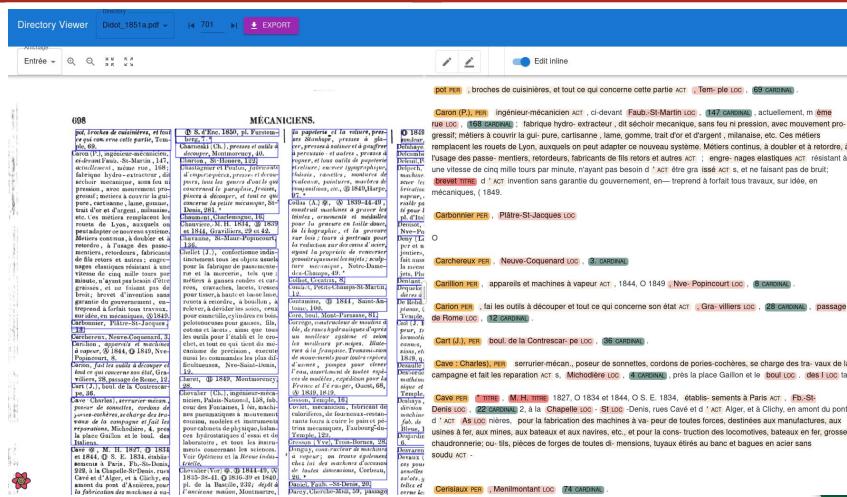
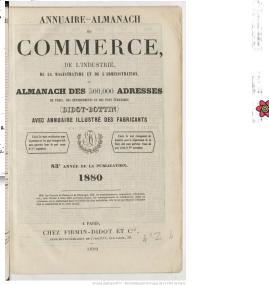
Explorer

Visualiser

Annoter

Analyser

API IIIF



Export structure

Assister le travail humain avec des outils automatiques

Explorer 

Visualiser 

Annoter 

Analyser 

API IIIF
E.g {BNF Gallica}

Export structuré
JSON / TEI XML

Interopérabilité

Des micro-services intégrés à l'annotation semi-automatique

TraITEMENT
automatique 

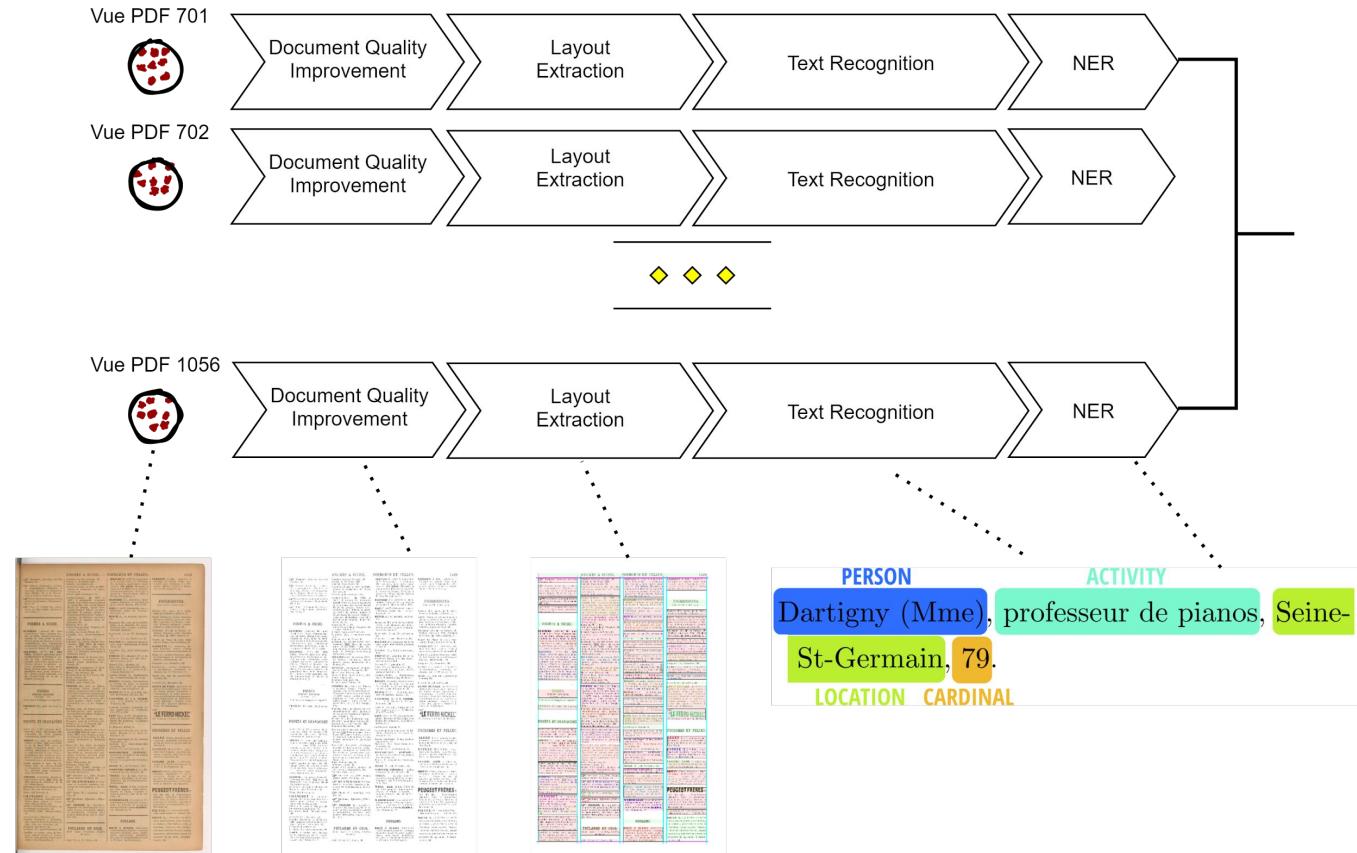


Correction & annotation
manuelle 

D'autres avantages:

- App web (distante ou locale) facile à déployer
- Des services de traitements externalisés (on peut changer et combiner des méthodes)

Chaîne de traitement automatique



Une reconnaissance perturbée par différentes imperfections



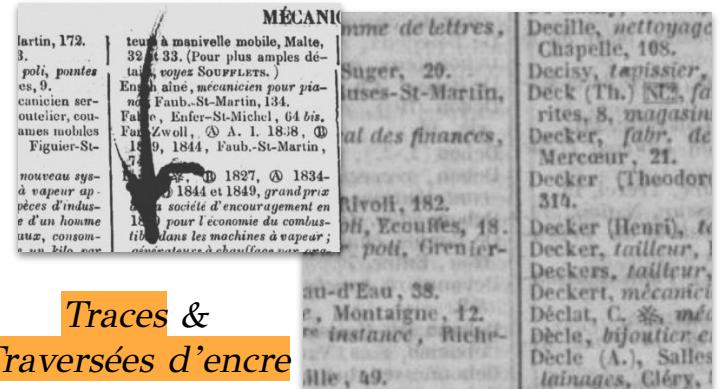
Archivage



Numérisation



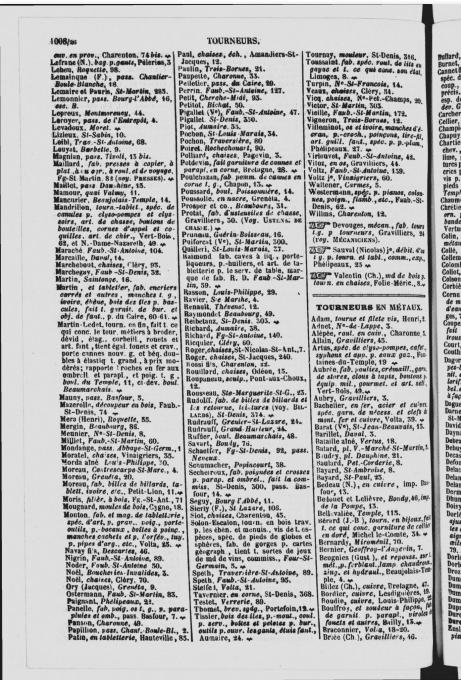
Compression



Traces & Traversées d'encre



Perte des niveaux de gris & Artefacts de compression



Inclinaison & Rotation & Courbure

Extraction de canvas / Segmentation

Méthode rapide basée XY-Cut pour les blocs

- Plusieurs niveaux hiérarchiques
 - Gestion efficace des colonnes/blocs qui exploite la redondance visuelle des pages

Méthode basée Watershed pour les lignes / entrées

- L'indentation des lignes / espaces de fin de lignes permettent de déterminer un début d'entrée

Límites

- Sensible au bruit, d'où le nettoyage en amont
 - Classification simultanée des régions
 - Information “logique” intégrée (gestion du multi-page...)

Reconnaissance du texte (OCR)

Mettereau, prop., quai d'Anjou, 7.
Mettemberg, élég., méd., St-Thomas-d'Enf., 5.
Metz (de), rentier, St-Guillaume, 30.
Metzinger, avocat, Rameau, 6.
Metzmacher, peint. sur émaux, St-Martin, 124.

Meurgey, épicier-herboriste, Dragon, 33.
Meurice, Chaussée-d'Antin, 3.
Meurice (Eug.), tapissier, Vivienne, 12.
Meurillon, marbrier-sculpteur, butte Mont-Parnasse, 15.



Analyse de la mise en page

Mettereau, prop., quai d'Anjou, 7.
Mettemberg, élég., méd., St-Thomas-d'Enf., 5.
Metz (de), rentier, St-Guillaume, 30.
Metzinger, avocat, Rameau, 6.
Metzmacher, peint. sur émaux, St-Martin, 124.

Meurgey, épicier-herboriste, Dragon, 33.
Meurice, Chaussée-d'Antin, 3.
Meurice (Eug.), tapissier, Vivienne, 12.
Meurillon, marbrier-sculpteur, butte Mont-Parnasse, 15.

OCR



Mettereau, prop., quai d'Anjou, 7.\nMettemberg, élég., méd., St-Thomas-d'Enf., 5.\nMetz (de), rentier, St-Guillaume, 30.\nMetzinger, avocat, Rameau, 6.\nMetzmacher, peint. sur émaux, St-Martin, 124.\nMeurgey, épicier-herboriste, Dragon, 33.\nMeurice, Chaussée-d'Antin, 3.\nMeurice (Eug.), tapissier, Vivienne, 12.\nMeurillon, marbrier-sculpteur, butte Mont-\nParnasse, 15.\n

La modularité de l'app nous a permis de tester plusieurs OCR rapidement (3 sont déjà API-fiés) :

- Kraken
- Tesseract ★★
- Pero OCR ★★

Reconnaître et structurer les informations des entrées

Composition des entrées:

- **Nom de la (des) personne(s) qui exerce(nt) une activité ou tien(nen)t un commerce.**
- Titre honorifique ou professionnel,
- Type d'activité ou de commerce,
- Type de local *,
- **Nom de voie *,**
- **Numéro de voie *,**
- Nom de section *.

L'ordre des informations peut varier selon le type d'index ou l'éditeur.

Seules les informations en gras sont systématiquement présentes.

Les informations suivies d'une * peuvent apparaître plusieurs fois par entrée.

Les étapes précédentes peuvent avoir produit un texte bruité : entrées mal délimitées, caractères mal reconnus, manquants, etc.



Une approche de reconnaissance des différents types d'informations (Named Entity Recognition) à base de règles serait coûteuse à développer et aurait peu de chances de traiter tous les cas de figure avec succès.

Reconnaissance d'entités nommées à base de réseaux de neurones profonds

1. Quels sont les modèles **récents disponibles** pour réaliser cette tâche de NER ?
2. Ces modèles de NER ont-ils **besoin de beaucoup de données d'entraînement** pour bien fonctionner sur les annuaires ?
3. Ces modèles de NER peuvent-ils produire de **bons résultats malgré le bruit OCR** ?
4. Peut-on **améliorer la résistance de ces modèles de NER au bruit OCR** ?

N. Abadie, E. Carlinet, Joseph Chazalon, B. Duméniel. A Benchmark of Named Entity Recognition Approaches in Historical Documents: Application to 19th Century French Directories. Document Analysis Systems. DAS 2022., May 2022, La Rochelle, France.

<https://github.com/soduco/paper-ner-bench-das22>

Modèles NER évalués

Nécessaire pour traiter des entités non standards

Entraînement original

modèle sur étagère

SpaCy CNN:

Pré-entraînement non supervisé: deep-sequoia Entraînement supervisé pour le NER: wikiner-fr

Adaptation de domaine

entraînement spécifique au corpus

Pré entraînement

Entraînement supervisé pour le NER: FTD labelled

CamemBERT:

Pré-entraînement non supervisé: OSCAR Entraînement supervisé pour le NER: wikiner-fr

Pré entraînement

Entraînement supervisé pour le NER: FTD labelled

CamemBERT pré-entraîné:

Pré-entraînement non supervisé: OSCAR Entraînement supervisé pour le NER: wikiner-fr

Pré-entraînement non supervisé:

FTD unlabelled

Entraînement supervisé pour le NER: FTD labelled

Les modèles testés sont disponibles sur nos dépôts HuggingFace et Zenodo ([10.5281/zenodo.6576008](https://doi.org/10.5281/zenodo.6576008)).

Transcriptions Pero OCR non corrigées pour 845,000 entrées brutes (≈ 7000 p.)

Transcriptions et annotations manuelles pour 8765 entrées (78 p.)

Quantité d'annotations et qualité des résultats

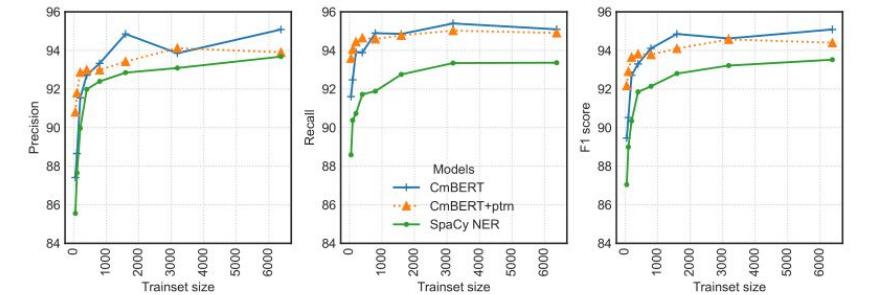
2. A-t-on besoin de beaucoup de données d'entraînement ?

Non ! On obtient de **bon résultats avec peu de données d'entraînement**, surtout si on réalise aussi un **pré-entraînement non supervisé**.

3. Les modèles de NER à base de réseaux de neurones profonds sont-ils adaptés pour traiter des annuaires hétérogènes ?

Oui ! Et les **modèles de type Transformer** obtiennent de très bon résultats.

	Training examples	49	99	199	398	796	1593	3186	6373
%		0.8	1.6	3.1	6.2	12.5	25.0	50.0	100.0
F1 score	CmBERT	89.5	90.5	92.7	93.3	94.1	94.9	94.6	95.1
	CmBERT+ptrn	92.2	92.9	93.6	93.8	93.8	94.1	94.6	94.4
	SpaCy NER	87.0	89.0	90.3	91.9	92.1	92.8	93.2	93.5



F1-scores moyens obtenus pour 5 entraînements + tests

Stratégie d'entraînement pour réduire l'influence du bruit OCR sur les résultats

Modèles testés :

- CamemBERT
- CamemBERT pré-entraîné

Entrées :

- Extraits d'annuaires corrigés et annotés manuellement
- Extraits d'annuaires directement produits par OCR (Tesseract V4 et Pero)

Sorties — Prédictions :

- prédictions NER de camemBERT
Avec et sans pré-entraînement

Sorties — Référence :

- Extraits d'annuaires corrigés et annotés manuellement
- Ou annotations manuelles projetées sur les extraits d'annuaires directement produits par OCR (bruités).

Mesures :

Jeux de données utilisés pour entraîner et tester CamemBERT ($\times 12$ variantes)

Jeux de pré-entraînement	Jeux d'entraînement	Jeux de test
<ul style="list-style-type: none">• Aucun• PERO (brut)	<ul style="list-style-type: none">• Référence• PERO*(* annotations projetées)	<ul style="list-style-type: none">• Référence• PERO*• Tesseract*(*annotations projetées)

Extrait du jeu de test de référence:

Annotations manuelles :

Dulay **PER**, chaudronnier **ACT**, r. du Pont- aux Choux **LOC**, 15
CARDINAL, 314.

Extrait du jeu de test avec annotations projetées sur du texte bruité :

Dulay **PER**, chandronnier **ACT**, +. du Pont-anx-Cars **LOC** Ge 7
CARDINAL Fe ÊR one

Stratégie d'entraînement pour réduire l'influence du bruit OCR sur les résultats

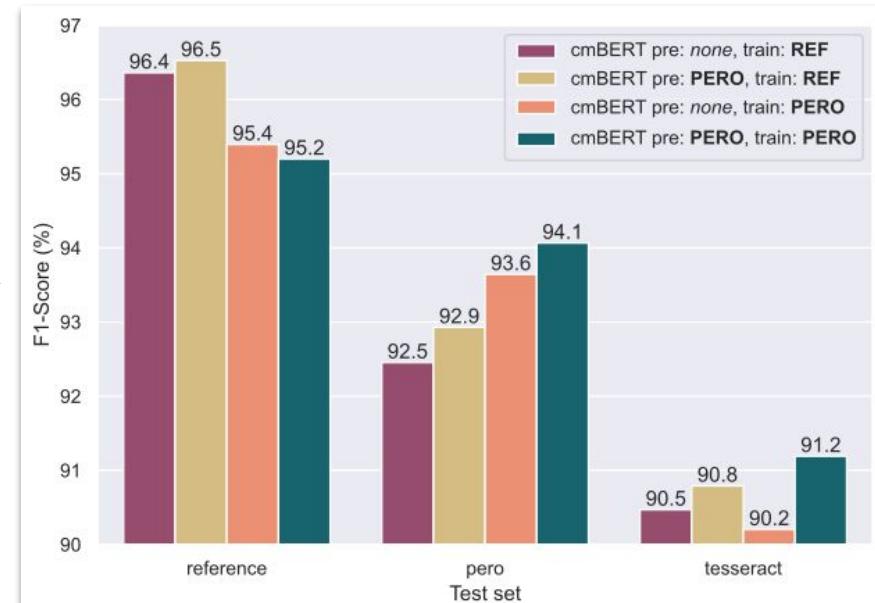
3 bis. Les modèles de NER à base de réseaux de neurones profonds sont-ils adaptés pour traiter des **transcriptions OCR bruitées** ?

Oui ! Mais les résultats sont moins bons que quand on a des textes propres.

4. Peut-on rendre les modèles de NER à base de réseaux de neurones profonds **plus robustes au bruit OCR** ?

Oui ! Il faut les pré-entraîner et les entraîner sur des textes bruités.

⇒ N'entraînez pas votre modèle sur du texte propre si vous voulez traiter des textes OC Risés !



F1-scores moyens obtenus pour 5 entraînements + tests

Démonstration

Quelques résultats qualitatifs

1. Cas positifs (ancienne interface)

View 700 Didot 1856a Force compute

Boulon, Charenton, 81.

Bouquet, M. H. Exposition universelle de 1855, breveté s. g. d. g., spécialité de tables de salles à mangèr, nouveau système mécanique donnant la facilité à une seule personne de pouvoir l'ouvrir, n'importe sa grandeur, Faub.-St-Antoine, 99, passage du Bras-d'Or. *

Bourdier, fab. de divans en tous genres, fait tout ce qui concerne sa partie, Charenton, 16.

Bourdier ainé, Nve-de-Lappe, 16.

Bourgade. Rumsford. 14.

Handles size

(0/159 entries checked)

ENTRY

OCR NER Form JSON

PER ACT LOC CARDINAL FT TITRE

CLEAR

Boulon, Charenton, 81.

Checked

Comment

• EAUX AUROPHILE ET ARGENTO- PHILE Berger, eaux spéciales pour nettoyages, déposées conformément à la loi. L'Aurophile, liquide pour nettoyer avec la plus grande facilité, propreté et sans altération, toutes pièces en bronze dorés, mates ou polies, montées ou non, avec marbre, etc., comme pendules, lustres, candélabres, garnitures de meubles, les statuettes en biscuit, marbre, etc., les fleurs ou autres sujets très-compliqués en porcelaine, les camées, livoire, la nacre, les cristaux polis ou dépolis, et avec beaucoup de précaution, les dorures sur bois. L'Argentophile, liquide pour nettoyer avec la plus grande facilité et sans altération toutes pièces d'orfèvrerie mates ou brunes, soit argent, argentées ou plaquées : la joaillerie, la bijouterie, ainsi que

View 700 Didot 1856a Local Server Force compute

Handles size

Bouquet, M. H. Exposition universelle de 1855, breveté s. g. d. g., spécialité de tables de salles à mangèr, nouveau système mécanique donnant la facilité à une seule personne de pouvoir l'ouvrir, n'importe sa grandeur, Faub.-St-Antoine, 99, passage du Bras-d'Or. *

Bourdier, fab. de divans en tous genres, fait tout ce qui concerne sa partie, Charenton, 16.

Bourdier ainé, Nve-de-Lappe, 16.

Bourgade, Rumford, 14.

Bourgoin. St-Antoine. 176.

(0/159 entries checked)

ENTRY

OCR NER Form JSON

CLEAR

Bouquet, M. H. Exposition universelle de 1855, breveté s. g. d. g., spécialité de tables de salles à mangèr, nouveau système mécanique donnant la facilité à une seule personne de pouvoir l'ouvrir, n'importe sa grandeur, Faub.-St-Antoine, 99, passage

 Checked

Comment

- EAUX AUROPHILE ET ARGENTO- PHILE Berger, eaux spéciales pour nettoyages, déposées conformément à la loi. L'Aurophile, liquide pour nettoyer avec la plus grande facilité, propreté et sans altération, toutes pièces en bronze dorés, mates ou polies, montées ou non, avec marbre, etc., comme pendules, lustres, candélabres, garnitures de meubles, les statuettes en biscuit, marbre, etc., les fleurs ou autres sujets très-compliqués en porcelaine, les camées, l'ivoire, la nacre, les cristaux polis ou dépolis, et avec beaucoup de précaution, les dorures sur bois. L'Argentophile, liquide pour nettoyer avec la plus grande facilité et sans altération toutes pièces d'orfèvrerie mates ou brunes, soit argent, argentées ou plaquées : la joaillerie, la bijouterie, ainsi que

View 700 Didot 1856a Local Server Force compute

Handles size

Bourdier, fab. de divans en tous genres, fait tout ce qui concerne sa partie, Charenton, 16.

Bourdier ainé, Nve-de-Lappe, 16.

Bourgade, Rumford, 14.

Bourgoin, St-Antoine, 176.

Bourlier, Charonne, 99.

Boutard, Traversière, 76.

Boutung (seule maison connue depuis trente ans), fab. et magasin de meubles en tous genres, modernes, antiques et fantaisie ; tient aussi le siège ; commission et exportation ; Faub.-St-Antoine, 97 - ci-devant même rue. 23. *

(0/159 entries checked)

ENTRY

OCR NER Form JSON

PER ACT LOC CARDINAL FT TITRE

CLEAR

Bourdier, fab. de divans en tous genres, fait tout ce qui concerne sa partie, Charenton, 16.

 Checked

Comment

- EAUX AUROPHILE ET ARGENTO-PHILE Berger, eaux spéciales pour nettoyages, déposées conformément à la loi. L'Aurophile, liquide pour nettoyer avec la plus grande facilité, propreté et sans altération, toutes pièces en bronze dorés, mates ou polies, montées ou non, avec marbre, etc., comme pendules, lustres, candélabres, garnitures de meubles, les statuettes en biscuit, marbre, etc., les fleurs ou autres sujets très-compliqués en porcelaine, les camées, l'ivoire, la nacre, les cristaux polis ou dépolis, et avec beaucoup de précaution, les dorures sur bois. L'Argentophile, liquide pour nettoyer avec la plus grande facilité et sans altération toutes pièces d'orfèvrerie mates ou brunes, soit argent, argentées ou plaquées : la joaillerie, la bijouterie, ainsi que

View 700 Didot 1856a Local Server Force compute

Bourdier aîné, Nve-de-Lappe, 16.

Bourgade, Rumford, 14.

Bourgoïn, St-Antoine, 176.

Bourlier, Charonne, 99.

Boutard, Traversière, 76.

Boutung (seule maison connue depuis trente ans), *fab. et magasin de meubles en tous genres, modernes, antiques et fantaisie ; tient aussi le siège ; commission et exportation* ; Faub.-St-Antoine, 97 ; ci-devant même rue, 23. *

Handles size

(0/159 entries checked)

ENTRY

OCR NER Form JSON

PER ACT LOC CARDINAL FT TITRE

CLEAR

Bourdier aîné, Nve-de-Lappe, 16.

Checked

Comment

- EAUX AUROPHILE ET ARGENTO- PHILE Berger, eaux spéciales pour nettoyages, déposées conformément à la loi. L'Aurophile, liquide pour nettoyer avec la plus grande facilité, propreté et sans altération, toutes pièces en bronze dorés, mates ou polies, montées ou non, avec marbre, etc., comme pendules, lustres, candélabres, garnitures de meubles, les statuettes en biscuit, marbre, etc., les fleurs ou autres sujets très-compliqués en porcelaine, les camées, l'ivoire, la nacre, les cristaux polis ou dépolis, et avec beaucoup de précaution, les dorures sur bois. L'Argentophile, liquide pour nettoyer avec la plus grande facilité et sans altération toutes pièces d'orfèvrerie mates ou bruniées, soit argent, argentées ou plaquées : la joaillerie, la bijouterie, ainsi que

View 700 Didot 1856a Force compute

Bourgade, Rumford, 14.
 Bourgoin, St-Antoine, 176.
 Bourlier, Charonne, 99.
 Boutard, Traversière, 76.
 Boutung (seule maison connue depuis trente ans), fab. et magasin de meubles en tous genres, modernes, antiques et fantaisie ; tient aussi le siège ; commission et exportation ; Faub.-St-Antoine, 97 ; ci-devant même rue, 23. *

Handles size

(0/159 entries checked)

ENTRY

OCR NER Form JSON

 PER ACT LOC CARDINAL FT TITRE

Bourgade, Rumford, 14.

 Checked

Comment

- EAUX AUROPHILE ET ARGENTO- PHILE Berger, eaux spéciales pour nettoyages, déposées conformément à la loi. L'Aurophile, liquide pour nettoyer avec la plus grande facilité, propreté et sans altération, toutes pièces en bronze dorés, mates ou polies, montées ou non, avec marbre, etc., comme pendules, lustres, candélabres, garnitures de meubles, les statuettes en biscuit, marbre, etc., les fleurs ou autres sujets très-compliqués en porcelaine, les camées, l'ivoire, la nacre, les cristaux polis ou dépolis, et avec beaucoup de précaution, les dorures sur bois. L'Argentophile, liquide pour nettoyer avec la plus grande facilité et sans altération toutes pièces d'orfèvrerie mates ou bruniées, soit argent, argentées ou plaquées : la joaillerie, la bijouterie, ainsi que

View 700 Didot 1856a Local Server Force compute

Bourgoin, St-Antoine, 176.

Bourlier, Charonne, 99.

Boutard, Traversière, 76.

Boutung (seule maison connue depuis trente ans), *fab. et magasin de meubles en tous genres, modernes, antiques et fantaisie ; tient aussi le siège ; commission et exportation* ; Faub.-St-Antoine, 97 ; ci-devant même rue, 23. *

Handles size

(0/159 entries checked)

ENTRY

OCR NER Form JSON

PER ACT LOC CARDINAL FT TITRE

CLEAR

Bourgoin, St-Antoine, 176.

Checked

Comment

- EAUX AUROPHILE ET ARGENTO- PHILE Berger, eaux spéciales pour nettoyages, déposées conformément à la loi. L'Aurophile, liquide pour nettoyer avec la plus grande facilité, propreté et sans altération, toutes pièces en bronze dorés, mates ou polies, montées ou non, avec marbre, etc., comme pendules, lustres, candélabres, garnitures de meubles, les statuettes en biscuit, marbre, etc., les fleurs ou autres sujets très-compliqués en porcelaine, les camées, l'ivoire, la nacre, les cristaux polis ou dépolis, et avec beaucoup de précaution, les dorures sur bois. L'Argentophile, liquide pour nettoyer avec la plus grande facilité et sans altération toutes pièces d'orfèvrerie mates ou bruniées, soit argent, argentées ou plaquées : la joaillerie, la bijouterie, ainsi que

Quelques résultats qualitatifs

2. Cas limites (nouvelle interface)

Figueran, Enghien, 7

Figueran PER , Enghien LOC , 7 CARDINAL .

Cas typique : entrée courte avec OCR et NER fonctionnels.

Poisat oncle et Cie , médailles
Ⓐ , fabrique d'acide sulfurique,
acide stéarique et oléique, et au-
tres produits, à la Folie-Nan-
terre, Enghien, 19. *

Poisat oncle et Cie PER , médailles , fabrique d'acide sulfurique, acide
stéarique et oléique, et au- tres produits ACT , à la Folie-Nan- LOC

* ACT terre, Enghien LOC , 19 CARDINAL . LOC

Entrée plus longue, avec un échec OCR (position "") qui entraîne un échec NER

Didot 1851a, vue 700

Berendorff(J.), M. H. 1844, Mouf-fetard, 294, atelier de construction de machines à vapeur et autres, (A) S. E., (A) 1844, (O) 1849.

Berendorff(J.) PER M. H. TITRE 18 44 TITRE , Mouf-fetard LOC , 294 CARDINAL , atelier de construction de machines à vapeur et autres ACT , S. TITRE E., O 1844, O 1849.

Cas typique : entrée longue avec entités principales (PER, LOC, CARD.) bien reconnues malgré l'inclusion de (","), et des erreurs NER sur les autres entités.

Duverneuil et La Tyna, 1806

Benaed , R. des Bons-Enfans , 19.

Benaed PER , R. des Bons-Enfans LOC , 19 CARDINAL .

Erreur OCR sur le nom ⇒ cible pour la correction par redondance

**Chaussard , R. de Grenelle, Halle-aux
Blés , 33.**

Chaussard PER , R. de Grenelle, Halle-aux Blés LOC , 33 CARDINAL .

Exemple positif, mais fusion de plusieurs niveaux de localisation

Deflandre 1829, vue 500

NOEL ainé et fils. Maison de commerce pour la
lithographie; r. N.-des-Petits-Champs , galerie
Colbert, esc. A.

NOEL ainé et fils PER . Maison de commerce pour la
lithographie ACT ; r. N.-des-Petits-Champs , galerie Colbert LOC . esc. A

Exemple positif, mais fusion de plusieurs niveaux de localisation

Favre et Duchesne, 1897

D'Huez, rue des Poulies, n.^o 209, —des Gardes-Françaises.

D'Huez PER , rue des Poulies LOC , n. 209 CARDINAL , — des
Gardes-Françaises LOC .

Micro-variations OCR sur la localisation (“—”), localisations multiples au même plan

des Commerçans de Paris.		49
Bourdon, <i>bouquetier</i> , r. S.-Denis, 269. (Eig.).	292	Bourgeois-Borlot, <i>plaquer</i> ; r. du Bac, 98.
Bourdon, <i>md. de beurre</i> , r. du Faub.-S.-Honoré, 46.	318	Bourgeois-Dumoulin, <i>md. de nouveautés</i> , r. Bussi, 2.
Bourdon et co., <i>commiss. par eau</i> , q. de la Mégiserie, 38.	325	Bourgeot, <i>fab. de chocolat</i> , r. S.-Honoré, 110.
Bourdon, <i>épicier</i> , r. des Quatre-Vents, 11.	351	Bourgeot, <i>papetier</i> , r. des Fossés-Montmartre, 31.
Bourdon (V.º), <i>épicière</i> , r. Sainte-Anne, 19.	351	Bourget et co., <i>commiss. de roulage</i> , r. S.-Denis, 152.
Bourdon (N.) alné et co., <i>miroitiers</i> , r. Bourg-l'Abbé, 48.	416	Bourget jeune, <i>commis. de roulage</i> , r. Beau- repaire, 3.
Bourdon (T.) <i>maîtrier</i> , r. Bourg-l'Abbé.		Bourrie, <i>épicier</i> , r. du Faub.-S.-Honoré

OCR

NER

318 Honoré, 46. Bourdon et co. PER , commiss. par cau ACT , f. de la
325 Megisserie LOC , 38 CARDINAL . Bussi, 2. 410 Bourge PER ot, fab.
de chocolat ACT , r. S.-Honoré LOC , 110 CARDINAL . 316

La détection de la structure reste une difficulté majeure.

Bilan

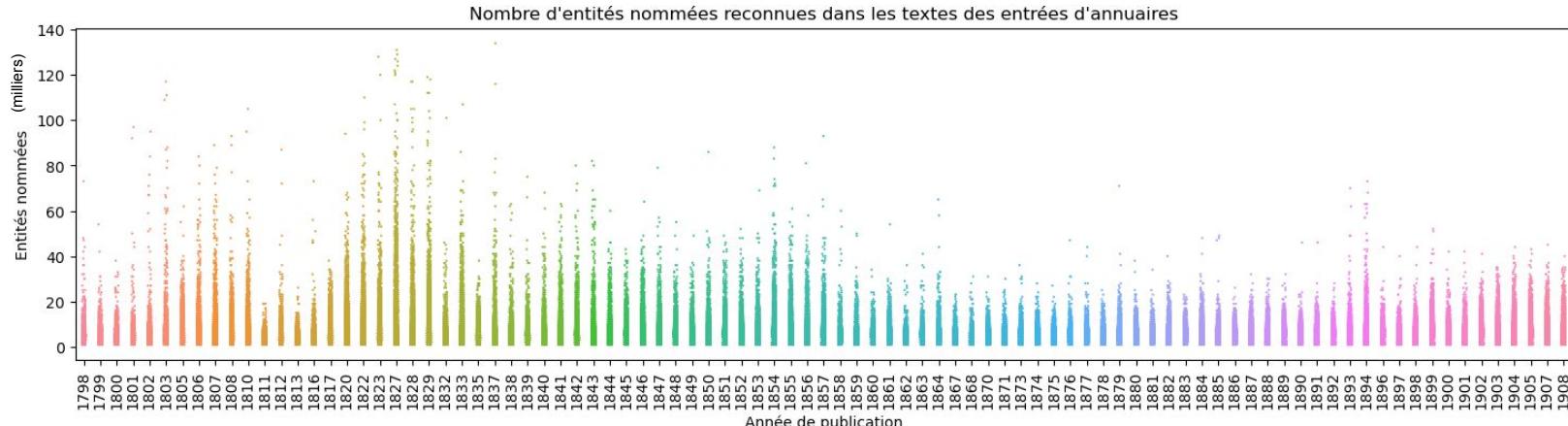
Vue d'ensemble des données extraites

Pour l'ensemble du corpus :

- 9 821 898 entrées extraites
- 7 260 104 (74%) entrées exploitables, de la forme : **PERSON** [**ACT**] **LOC** [**CARDINAL**]

⇒ la forte redondance des entrées est riche en information

Un enjeu restant : identifier, localiser et re-traiter (semi-automatiquement ?) les entrées mal identifiées



“Traîne” des entrées aggrégées : une faible proportion mais des effets potentiellement importants sur les traitements suivants.

Perspectives — Projet Biblissima+

Encore une plateforme d'annotation / reconnaissance de documents imprimés ?

⇒ Non ! Vers une **collection** d'outils **ouverts, minimaux**, laissant les historiens **constituer leurs corpus** et y “superposer” **leurs données**.

Nos objectifs à venir :

- **Consolider notre interface** de visualisation et d'annotation (page web unique, sans besoins d'hébergement ni d'installation)
- **Déployer des services** pour une utilisation à la demande sans nécessiter une machine puissante (OCR, NER, etc. — un hébergement HumaNum ?)
- **Enrichir notre boîte à outils** (correction de courbure des pages, export TEI...)
- **Partager** plus largement ces outils (*open source*, documentation, voire formation)