

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



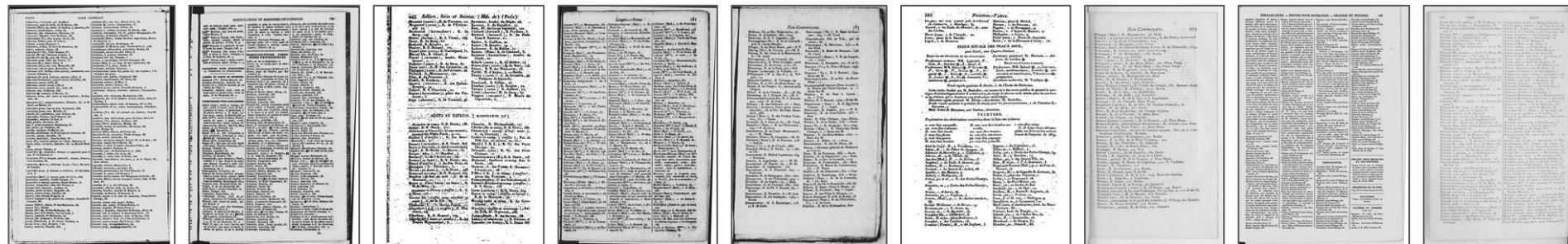
Introduction

Un corpus de **141** annuaires numérisés :

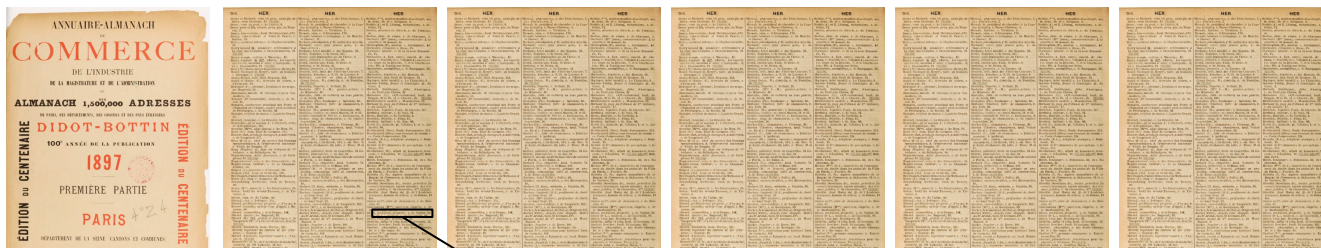
- Plusieurs centaines de milliers de pages,
- Plusieurs millions d'**entrées**, comportant le nom des commerces, leur type d'activité et leur localisation.

⇒ 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



Pages d'annuaires numérisées

Identification des entrées

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

Entrée d'annuaire

Transcription du texte

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

Texte numérique de l'entrée

Annotation du texte

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

Entités nommées

Structuration des entrées



Base de données intégrée

Construire un outil pour les historiens

Explorer



Visualiser



Annoter

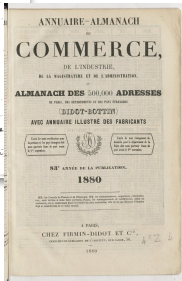


Analyser



API IIF

E.g {BNF Gallica



The screenshot shows a digital document viewer interface. At the top, there's a 'Directory Viewer' header with a file name 'Doc_1881.pdf', a page number '14 / 791', and an 'EXPORT' button. Below the header, there's a search bar with 'Entrée' and several icons. The main content area displays a page from the 'ALMANACH DU COMMERCE' with a section titled 'MÉCANICIENS'. The text on the page is annotated with red and blue markers. To the right of the document, there's a sidebar containing a list of links, each with a colored dot and text such as 'Carboneur act.', 'Cachoux res.', 'Cannon res.', etc. The interface also includes a 'Edit inline' button and a search bar at the bottom.



Export structuré

{JSON}
TEI <XML />



Assister le travail humain avec des outils automatiques

Explorer



Visualiser



Annoter



Analyser



API IIIF

E.g {BNF Gallica

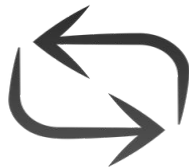
Export structuré

JSON / TEI XML

Interopérabilité

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

Traitement
automatique



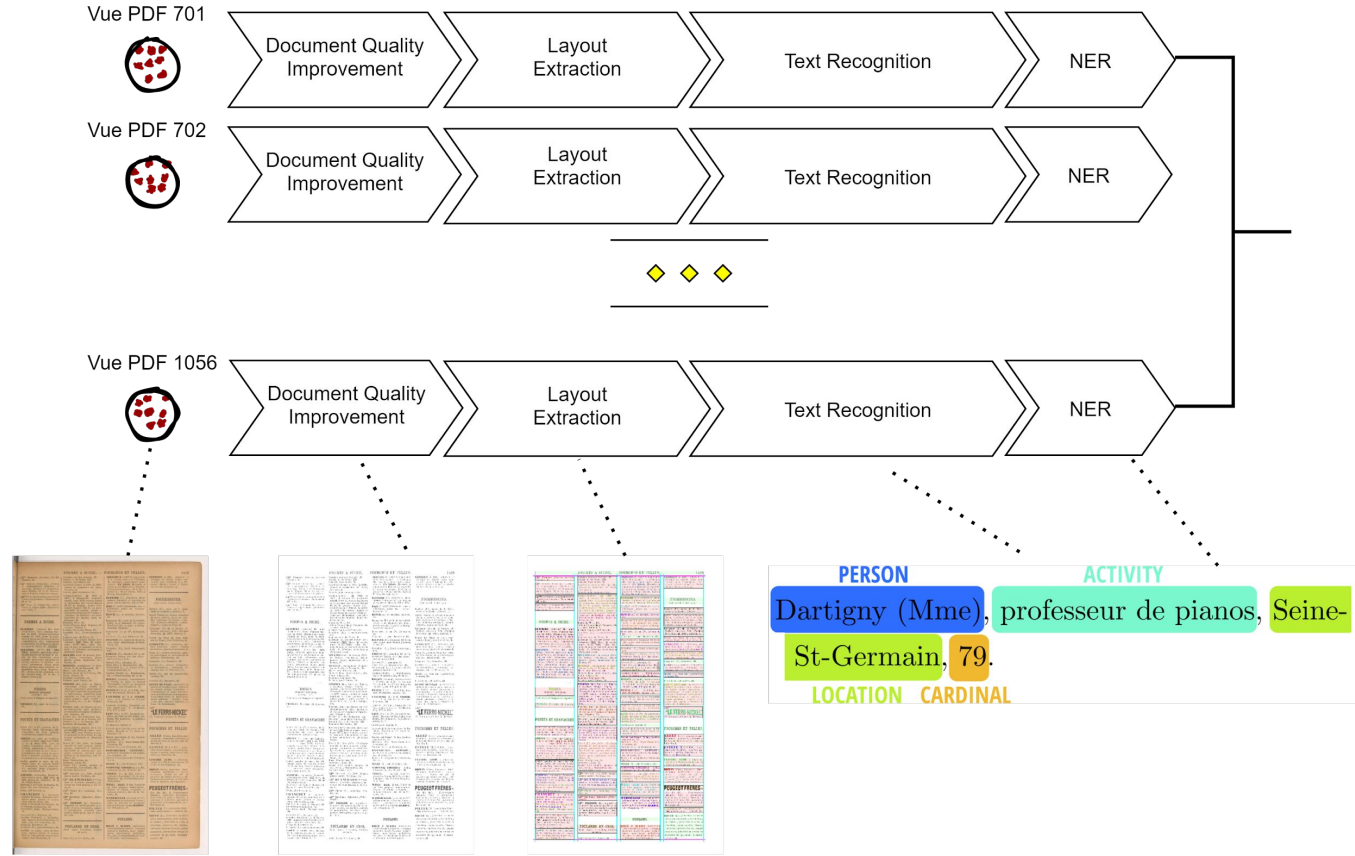
Correction & annotation
manuelle



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



MÉCANIQUE
... à manivelle mobile, Malte, 32 et 33. (Pour plus amples détails, voyez SOUFFLETS.)
Engh aîné, mécanicien pour pianicien, Faub.-St-Martin, 134.
Fauve, Eifer-St-Michel, 64 bis, Faub. Zwill, @ A. 1. 1838, @ 1839, 1844, Faub.-St-Martin, 187.
... 1827, @ 1834-1844 et 1849, grand prix de la société d'encouragement en 1849 pour l'économie de combustible dans les machines à vapeur; inventeur de chaudières sur eau.

Traces & Traversées d'encre

MERCIE RS.
Bollard, R. du Bouloy, 2.
Boniface, R. des Passes, S. Germain des Prés, 2.
Bouchard, R. du Roule, 16.
Boucher, (de noueautés) Boulev. du Panorama, 155.
Bouffin, R. aux Fers, 98.
Bourdin, (Mme), R. de la Loi, 24.
Bournet-Aubertot, R. des Moines, 22.
Boussange, R. des F. S. Germain l'Aux. 6.
Boutillier, (Mme) Pal. du Trib. Gal. de pierre, 155.
Boutier, R. S. Denis, 98.
Braschi, (noueautés) Pal. du Trib. Gal. de pierre, 103.
Bréant, R. Mâtinatoire, 84.

Decille, nettoyage
Chapelle, 108.
Decisy, tapisier,
Deck (Th.) St. fa
rites, 8, magasin
Decker, fabr. de
Mereaur, 21.
Decker (Theodor)
314.
Decker (Henri), t
Decker, tailleur,
Deckers, tailleur,
Deckert, mécanici
Déclat, C. mé
Déclé, bijoutier e
Declé (A.), Salles
lainages, Cléry,

Perte des niveaux de gris & Artefacts de compression

Inclinaison & Rotation & Courbure

TOURNEURS.
TOURNEURS EN MÉTAUX.
TOURNEURS EN BOIS.

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

Limites

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

FORMES A SUCRE.	FOURCHES ET PELLER.	1410
BIGNIER (R), succ. de E. Gagnan, 35. * Broomer, ferrière, rue des Gagnans, 35. * BIGNIER (R), succ. de E. Gagnan, 35. * BIGNIER (R), succ. de E. Gagnan, 35.	BOULET (R), 101, C. G. Lagomotti et C. Avinier, succ. de V. Vienneau, 20, caducière, montantes pour diam. EN GROS. Monnaie d'Argent, L. Marcoux et Zurich, 101, rue de Valenciennes, 11 et 15 et L. Mondier, Wood street, 123, Berlin, Hainlaube, 4.	BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris, succ. de J. Dupuy, 4, rue de Valenciennes, 11 et 15, Paris. BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris. BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris. BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris.
BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2.	BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris.	BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris.
BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2.	BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris.	BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris.
BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2.	BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris.	BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris.
BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2. * BOULET (R), succ. de F. Delecoeur, 2, rue de Valenciennes, 2.	BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris.	BOULET & Cie , appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris. * BOULET & Cie, appâts, 141, rue de Valenciennes, 11 et 15, Paris.

Reconnaissance du texte (OCR)

Mettereau, prop., quai d'Anjou, 7.	Meurgey, épicier-herboriste, Dragon, 33.
Mettemberg, élig., méd., St-Thomas-d'Enf., 5.	Meurice, Chaussée-d'Antin, 3.
Metz (de), rentier, St-Guillaume, 30.	Meurice (Eug.), tapissier, Vivienne, 12.
Metzinger, avocat, Rameau, 6.	Meurillon, marbrier-sculpteur, butte Mont-Parnasse, 15.
Metzmacher, peint. surémaux, St-Martin, 124.	

Analyse de la mise en page

Mettereau, prop., quai d'Anjou, 7.	Meurgey, épicier-herboriste, Dragon, 33.
Mettemberg, élig., méd., St-Thomas-d'Enf., 5.	Meurice, Chaussée-d'Antin, 3.
Metz (de), rentier, St-Guillaume, 30.	Meurice (Eug.), tapissier, Vivienne, 12.
Metzinger, avocat, Rameau, 6.	Meurillon, marbrier-sculpteur, butte Mont-Parnasse, 15.
Metzmacher, peint. surémaux, St-Martin, 124.	

OCR

Mettereau, prop., quai d'Anjou, 7.\nMettemberg, élig., méd., St-Thomas-d'Enf., 5.\nMetz (de), rentier, St-Guillaume, 30.\nMetzinger, avocat, Rameau, 6.\nMetzmacher, peint. surémaux, 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. Duménieu. 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 <i>modèle sur étagère</i>		Adaptation de domaine <i>entraînement spécifique au corpus</i>	
SpaCy CNN:	Pré-entraînement non supervisé: <u>deep-sequoia</u>	Entraînement supervisé pour le NER: <u>wikiner-fr</u>	Pré-entraînement	Entraînement supervisé pour le NER: <u>FTD labelled</u>
CamemBERT:	Pré-entraînement non supervisé: <u>OSCAR</u>	Entraînement supervisé pour le NER: <u>wikiner-fr</u>	Pré-entraînement	Entraînement supervisé pour le NER: <u>FTD labelled</u>
CamemBERT <i>pré-entraîné</i>:	Pré-entraînement non supervisé: <u>OSCAR</u>	Entraînement supervisé pour le NER: <u>wikiner-fr</u>	Pré-entraînement non supervisé: <u>FTD unlabelled</u>	Entraînement supervisé pour le NER: <u>FTD labelled</u>

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

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

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

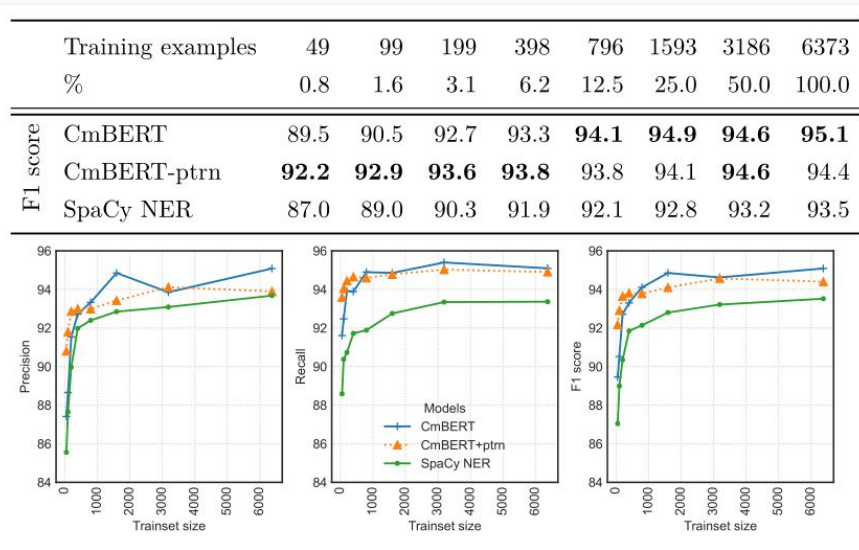
Quantité d'annotations et qualité des résultats

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.



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 : score F1

Jeux de données utilisés pour entraîner et tester CamemBERT (× 12 variantes)

Jeux de pré-entraînement	Jeux d'entraînement	Jeux de test
<ul style="list-style-type: none">• <i>Aucun</i>• 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

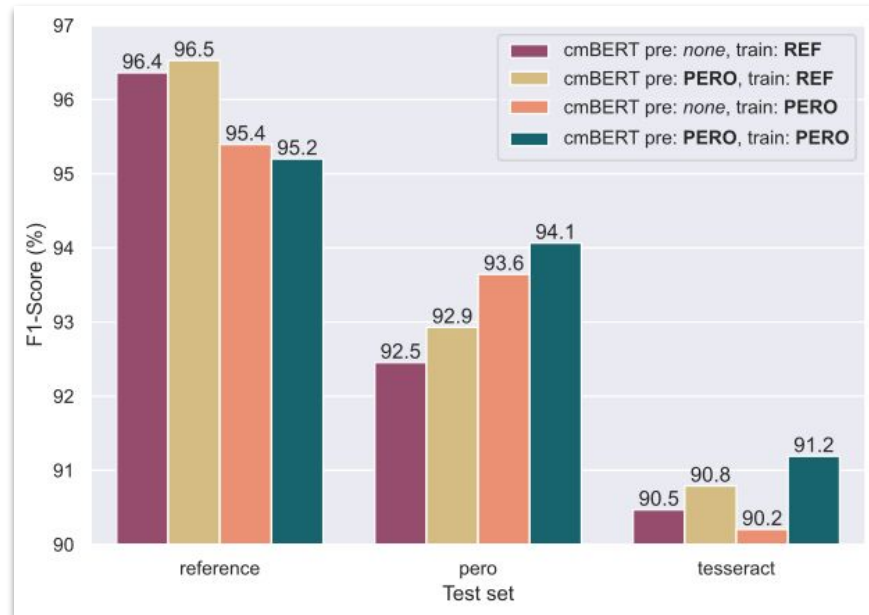
3 bis. Les modèles de NER à base de réseaux de neurones profonds sont-ils adaptés pour traiter des **transcriptions OCR bruitées** ?

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 OCRisés !



F1-scores moyens obtenus pour 5 entraînements + tests

Démonstration

Cabanis et Cie PER , association pour la construction des machines à éle- ver les eaux, pompes domestiques, moteurs hydrauliques de toute na- ture, machines soufflantes pour les forges, machines à eaux de Sellz, presses à découper, outillage et réparation de machines ACT , Vinai- griers LOC , 32 CARDINAL .

Comment
Tags

Checked

SAVE

MÉCANICIENS. 697

lons ronds et demi-ronds et carvés, et de toute sorte de moulure ; spécialité de fabrication de broches en acier et cuivre pour les peignes à tisser draps, voiles et ouvertures, toile métallique pour tout ce qui concerne le tissage, M. H. 1849, Amandiers-Popincourt, 19.

Auger (Vve), mécanicien, fabr. d'importe- pièces pour festons de garnitures de robes et de mantelets, importe- pièces pour cartonnage, semelles de souliers, etc., tient un assortiment de maillets, billets et plombs, pass. de la Trinité, 77.

Avoyne-Bainée, serrurier-mécan., fab. de lits en f. r., 1839-14, 1819, Boulangers-St-Victor, 22.

Bijard, quai Jemmapes, 248.

Baptrosse pere, tout ce qui a rapport à la fabrique d'indennes, Charenton, 58.

Barbot, Popincourt, 58.

Barbès, Trois-Bornes, 21.

Bardeu, Descartes, 38.

Barraud, horloger-mécanicien, entreprend toute espèce de fabrication nouvelle, quelle que soit la précision demandée, et construit toutes sortes de machines et outils pour Messieurs les inventeurs; vis cylindriques et coniques, découpage à façon, fendage de roues; pièces détachées fabriquées par procédés mécaniques pour toute espèce d'industrie, 1849, Ecoiffes, 20. au Mar., près le marché des Blancs-Manteaux. *

Baudat, constructeur de mécaniques à scier le bois en

Becker, mécanicien-graveur, breveté d'invention (sans garantie du gouvernement), fabrique spéciale de presses à copier, presses à timbre sec et bumide, presses à cacheter, presses à perçusion pour l'extraction des matières végétales et pour le satinage des papiers, composteurs pour papiers à lettres, pour raisons de commerce, griffes de toutes sortes, et généralement toute la gravure sur métaux, St-Denis, 380, passage Lemoine. *

Beerstecher (A.), Lourcine, 88.

Berendorff (J.), M. H. 1844, Moulletard, 294, atelier de construction de machines à vapeur et autres, 1849, S. E., 1844, 1849.

Bergeron, 1849, Marais-du-Temple, 73.

Bernier aîné, outils établis, presses, 1844, Faubourg-St-Antoine, 89.

Bernier, Menilmontant, 90.

Berthet, Simon-le-Franc, 13.

Berton, Neuve-St-Denis, 12.

Bertrand, mécanicien, breveté (sans garantie du gouvernement), construit toutes espèces de machines et outils, spécialité de machines pour les fabricants de matières premières de chapellerie, machines à vapeur et tours pour les fabricants de chapeaux; tondeuses et souffleuses de poils de lièvres et de lapins, réparation et entretien de machines à vapeur; on trouve toujours chez lui des machines à vapeur toutes prêtes à être mises en place, Vieille-du-Temple, 58, anc. 72

Boslay (Ch.) 1839 1849

provinc, à des prix très-moderés, St-Jacques, 261. *

Bosche aîné, ingénieur-mécanicien, inventeur de divers perfectionnements au métier Jacquard, S. E. 1850, Amandiers-Popincourt, 22, ci-devant St-Maur, 14. *

Boscher, Bondy, 70.

Bosquillon, 1823-27-34, constructeur de perçages accélérés et de mécaniques Jacquard parisiennes, Paradis-Poissonnière, 20.

Boucher, Saint-Pierre-Popincourt, 18.

Bouchon, dépôt de moulins à bras portatifs, Nve-St-Nicolas, 16.

Boulley (Ene), mécanicien, fab. de machines à couper les peaux de lapins, souffleuses de toutes dimensions marchant à bras et à la vapeur, tours de chapeliers, 1849, Francs-Bourgeois-Marais, 3, ci-devant Beaubourg, 59. *

Bouhon, appareil dit cale à léau, breveté sans garantie du gouvernement, destiné au soulagement des chevaux dans les montées, 1849, Société d'encouragement 1850, pl. Dauphine, 71. *

Bourdier, bandagiste-herniaire, mécanicien-orthopédiste, fournisseur de la Société protestante de prévoyance et de secours mutuels de Paris, fabrique spéciale de toutes espèces de ressorts et de bandages, béquilles, appareils orthopédiques, ceintures hypogastriques, bas lacés en tous genres, suspensoirs, objets d'allaitement, et tous articles en gomme élastique, Cadran, 5.

gnier, mécanique nouvelle, presse à copier les lettres, perfectionnée 1849, Saint-Jacques, 169, près le Panthéon. *

Brocard (Léon), machines à papier, St-Gilles, 12.

Burckel, fab. étrilles, fuseaux et broches, Nve-de-Lappe, 3.

Busser, fab. de presses lithographiques, machines à réparer, fait les réparations à des prix modérés, Lamartine, 18-20, ci-devant Coquenard. *

Butt, Buisson-St-Louis, 10.

Cabanis et Cie, association pour la construction des machines à élever les eaux, pompes domestiques, moteurs hydrauliques de toute nature, machines soufflantes pour les forges, machines à eau de Sellz, presses à découper, outillage et réparation de machines, Vinai-griers, 32.

Chouet, breveté (sans garantie du gouvern.) pour les porte-moules à bougies emboutis d'un seul morceau, et les moules montés sans aucune soudure, et pouvant les monter et les démonter soi-même; fabricant de mandrins, fondeur de moules de chandelles, bougies, cierges, spécialité, 1849, Marché-aux-Vaux, 4.

Calamel, Dupetit-Thouars, 10.

Calard (T.-F.), fabr. de feuilles métalliques percées en fer, cuivre et zinc; tôle en râpe pour le nettoyage des blés, cribles à main pour remplacer la peau.—Atelier de construction de machines à percer, à cintrer les fers à froid et à ajuster les ressorts de voitures. Leclerc. 8. Faub.-St-Jac-



Quelques résultats qualitatifs

1. Cas positifs (ancienne interface)

View 700 Didot 1856a Local Server Handles size Force compute

Boulon, Charenton, 81.

Bouquet, M. H. Exposition universelle de 1855, breveté s. g. d. g., spécialité de tables de salles à manger, 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 aîné, Nve-de-Lappe, 16.

Bourgade. Rumford. 14.

(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, 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 brunies, soit argent, argentées ou plaquées : la joaillerie, la bijouterie, ainsi que

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 aîné, Nve-de-Lappe, 16.

Bourgade, Rumford, 14.

Bourgoin. St-Antoine. 176.

(0/159 entries checked)

ENTRY

OCR NER Form JSON

PER

ACT

LOC

CARDINAL

FT

TITRE

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 dures 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 Handles size Force compute

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

Bourdier aîné, 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 dures 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

Handles size

Force compute

Bourdier aîné, Nve-de-Lappe, 16.

Bourgade, Rumford, 14.

Bourgoin, St-Antoine, 176.

Bourlier, Charonne, 99.

Boutard, Traversière, 76.

Boutang (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 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 brunies, soit argent, argentées ou plaquées : la joaillerie, la bijouterie, ainsi que

View 700 Didot 1856a Local Server Handles size Force compute

Bourgade, Rumford, 14.
 Bourgoin, St-Antoine, 176.
 Bourlier, Charonne, 99.
 Boutard, Traversière, 76.
 Boutang (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

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 brunies, soit argent, argentées ou plaquées : la joaillerie, la bijouterie, ainsi que

soduco-annot-gui

View 700 Didot 1856a Local Server Handles size Force compute

Bourgoin, St-Antoine, 176.
Bourlier, Charonne, 99.
Boutard, Traversière, 76.
Boutang (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

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 brunies, soit argent, argentées ou plaquées : la joaillerie, la bijouterie, ainsi que

Quelques résultats qualitatifs
2. Cas limites (nouvelle interface)

Didot 1851a, vue 800

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*,
Ⓐ S. E., Ⓐ 1844, Ⓞ 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 aîné et fils. Maison de commerce pour la lithographie; r. N.-des-Petits-Champs , galerie Colbert, esc. A.

NOEL aîné 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.° 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

Bottin, 1820, vue 52

<i>des Commerçans de Paris.</i>		47
Bourdon, <i>bonnetier</i> , r. S.-Denis, 269.	Bourgeois-Borlot, <i>plaqueur</i> , r. du Bac,	
(Elig.), 292	98.	451
Bourdon, <i>md. de beurre</i> , r. du Faub.-S.-	Bourgeois-Damoulin, <i>md. de nouveautés</i> , r.	
Honoré, 46.	Bussi, 2.	419
Bourdon et co., <i>commiss. par eau</i> , r. de la	Bourgeot, <i>fab. de chocolat</i> , r. S.-Honoré,	
Mégisserie, 38.	110.	316
Bourdon, <i>épiciier</i> , r. des Quatre-Vents,	Bourgeot, <i>papetier</i> , r. des Fossés-Montmar-	
11.	tre, 31.	437
Bourdon (V.), <i>épicière</i> , r. Sainte-Anne,	Bourget et co., <i>commiss. de roulage</i> , r. S.-	
10.	Denis, 152.	323, 381
Bourdon (N.) aîné et co., <i>mirroitiers</i> , r.	Bourget jeune, <i>commiss. de roulage</i> , r. Beau-	
Bourg-l'Abbé, 48.	repaire, 3.	323
Bourdon (N.) <i>quincaillier</i> , r. Bourg-l'Abbé.	Bourgeois, <i>épiciier</i> , r. du Faub.-S.-Honoré	

La détection de la structure reste une difficulté majeure.

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

Bilan

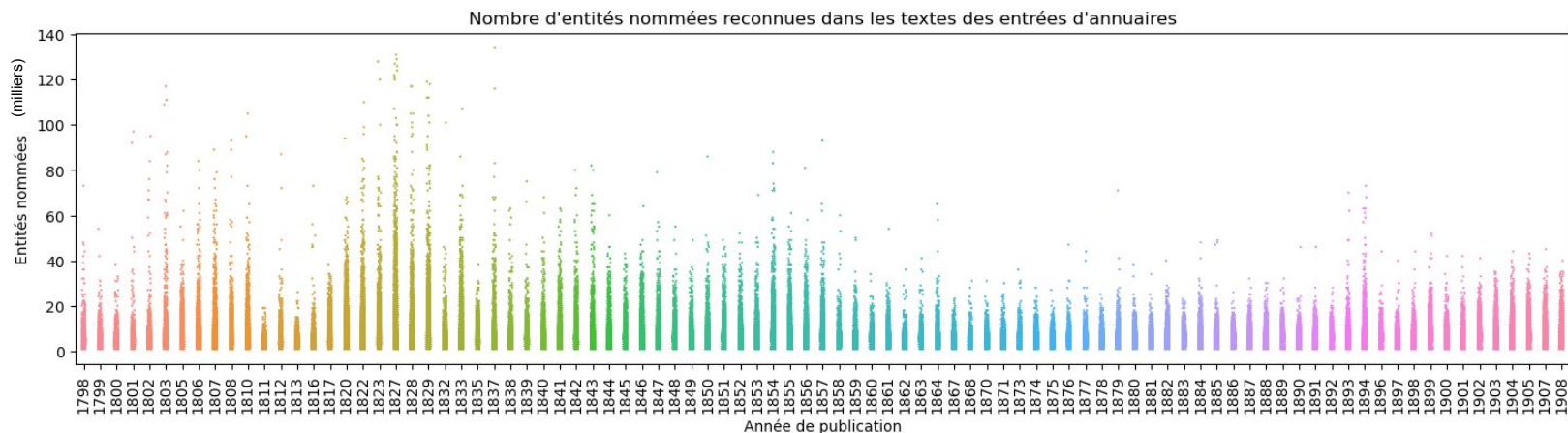
Vue d'ensemble des données extraites

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 agglomérées : 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)