

Extraction de données

Approche du projet SoDUCo

De l'image d'annuaire du commerce ancien
au graphe de données interrogeable

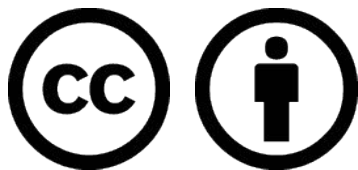
N. Abadie, S. Baciocchi, E. Carlinet, J. Chazalon, P. Cristofoli, B. Duménieu, J. Perret, S. Tual
Lastig (IGN), CRH (EHESS), LRE (EPITA)

6 nov. 2023, séminaire SoDUCo-BnF



(Ré)utiliser cette présentation

Licence **Creative Common BY 4.0**



N'hésitez pas à **réutiliser** ces contenus et nous **citer**.
*Sauf citation ou mention contraire,
nous sommes les auteurs des contenus.*

Consultez le diaporama pendant et après la présentation

<https://bit.ly/soduco-bnf-20231106-extraction>



Un corpus hétérogène

Ouvrages produits par des éditeurs variés, numérisés par différents prestataires.

Des différences de

- Longueur, richesse et structuration de l'information
- Mises en page, fontes
- Qualité inégale : papier, résolution des images, présence de marques...

Non-Commerçans. (Paris). 269

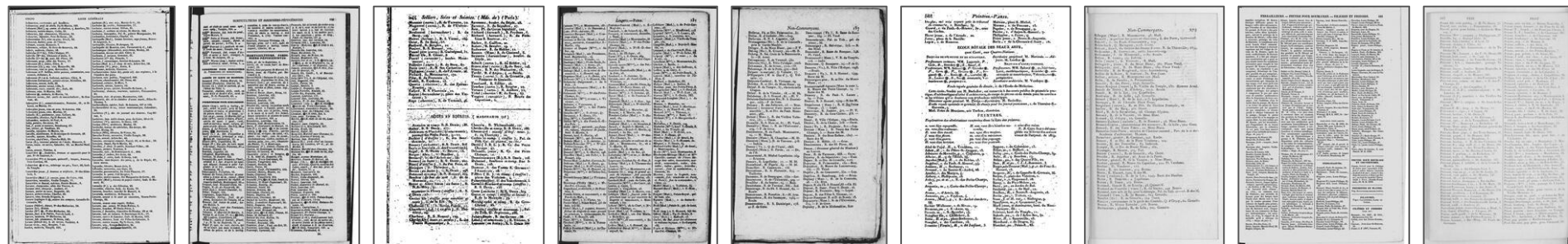
<p>Chardin, R. Pavée, 16. — R. G. Chardin, R. Michel Lepelletier, 21. Chardon, (Ve.) R. S. Marc, 15.</p>	<p>Cheviron, R. Chapou, 13. Chimay, (Mme.) R. de Varennes, 31. Choart-Duplessis, R. de Turanne, 31.</p>
---	--

AMADOU ET ALLUMETTES. — Pour les ALLUMETTES OXIGÉNÉES.
 Voyez BRIQUETS PHYSIQUES.

DARRAS (Thomas), r. de la Vieille-Monnaie, 10. GALLIENNE j^e, r. de la Heumerie, 3.

Briquets et veilleuses, mèches à quinquets, à quinquets, à quinquet, veilleuses mèches, souffrées ; mèches souffrées, pierres, agaric de chêne, pierres, agaric, bouchons, liège.
 Liège en planches, bouchons. LEROY, r. Aubry-le-Boucher, 43.

<p>BAUDOYER (place). IX Arr. Hôtel-de-Ville.) ← Rue Tixeranderie, pourtour St-Gervais, Saint-Antoine et Renaud-LeFèvre.</p>	<p>26* Longpré aîné, bijoutier en or et argent. Saint-Omer, émailleur. Cellier (A.), graveur-ciseleur.* Rousseau (J.), bijoutier en or.* Benoit, orfèvre-fabr. Léréys, doreur.</p>	<p>Bourguille, fabr. de presses. Vautain, passementier. Fininojac, bronze doré. Rabé aîné, fabr. de bal- lons.* Gaulin, chapelier. Moisy, tabletier.</p>	<p>7 Ecole communale de jeunes filles. Berthelot, vins. 6 Verstaen, serrurier-mécanicien. 8 Michel, brossier. 9 Labottiere, serrurier. 10 Sacrez, vins. 12 Baudoin, épici. 13 Lejard, clouteries et crêpines. 14 Badauel (Vve), fab. de</p>
<p>1 Lissoty (Vve), vins. 2* Privé, distillateur. Lemoine-Cluzel et Leroy, nouveautés. Chantrier, court-gourm.</p>	<p>30 Bouton, fab. de cuir vernis.* 31 Pardon, vins.</p>	<p>49* Cendrier aîné, prop. Desmarests, fab. bottes d'emballage. Ferrand, lapidaire.</p>	<p>et tapisseries. 10 Lainé jeune, vins. Jumelles omnibus et entrepise générale des Omnibus. 11 Melouzy, vins en gros, et à Bercy, Port, 31. 12 Combaud, coiffeur. Monmain (P.), vins en gros. 13 Dufailly, sculpt. fabr. de carton-pierre.</p>



Un instantané des acteurs commerciaux à un instant

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibl'que protestante (Société), Moulins, 16.
Bibron, aide-natural., au Muséum d'hist. nat.
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bicant (Mme), fondeur en cuivre, cour de la Corderie-du-Temple, 26.
Bichard (Mme), Nve-de-Luxembourg, 17.
Bichard, tabacs et eau-de-vie, Faub.-St-Martin, 45.

Didot 1841a - page 95

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibl'que protestante (Société), Moulins, 16.
Bibron, aide-natural., au Muséum d'hist. nat.
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bican (Vve) et fils, fondeur en cuivre, place de la Corderie-du-Temple, 26.
Bicel, épicier, marché d'Aguesseau, 15.
Bichard (Mme), Nve-de-Luxembourg, 17.
Bichard, tabacs et eau-de-vie, Faub.-St-Martin, 45.

Didot 1842a - page 117

Bibliothèque Ste-Geneviève, rue des Sept-Voies et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bical et Doire, fab. de socques, Vert-Bois, 14.
Bican (Vve) et fils, fondeurs en cuivre, place de la Corderie-du-Temple, 26.
Bicel, épicier, Marché-d'Aguesseau, 15.
Bichard, tabac et eau-de-vie, Faub.-St-Martin, 45.

Didot 1843a - page 129

Bibliothèque Ste-Geneviève, rue des Sept-Voies et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibolet, relieur, passage Sainte-Marie-Saint-Germain, 10.
Bibonne, architecte, Magasins, 12.
Bibron, aide-naturaliste au Jardin-des-Plantes, Cuvier, 29.
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bican (Vve) et fils, fondeurs en cuivre, place de la Corderie-du-Temple, 26.
Bicel, épicier, Marché-d'Aguesseau, 15.
Bichard, tabac et eau-de-vie, Faub.-St-Martin, 45.

Didot 1844a - pages 125-126

Un instantané des acteurs commerciaux à un instant

Grande **redondance** d'information intra- et inter-annuaires.

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibl'que protestante (Société), Moulins, 16.
Bibron, aide-natural., au Muséum d'hist. nat.
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bican (Mme), fondeur en cuivre, cour de la Corderie-du-Temple, 26.
Richard (Mme), Nve-de-Luxembourg, 17.
Bichard, tabacs et eau-de-vie, Faub.-St-Martin, 45.

Didot 1841a - page 95

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bib'que protestante (Société), Moulins, 16.
Bibron, aide-natural., au Muséum d'hist. nat.
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bican (Vve) et fils, fondeur en cuivre, place de la Corderie-du-Temple, 26.
Bicel, épicier, marché d'Aguesseau, 15.
Richard (Mme), Nve-de-Luxembourg, 17.
Bichard, tabacs et eau-de-vie, Faub.-St-Martin, 45.

Didot 1842a - page 117

Bibliothèque Ste-Geneviève, rue des Sept-Voies et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bical et Doire, fab. de socques, Vert-Bois, 14.
Bican (Vve) et fils, fondeurs en cuivre, place de la Corderie-du-Temple, 26.
Bicel, épicier, Marché-d'Aguesseau, 15.
Bichard, tabac et eau-de-vie, Faub.-St-Martin, 45.

Didot 1843a - page 129

Bibliothèque Ste-Geneviève, rue des Sept-Voies et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibolet, relieur, passage Sainte-Marie-Saint-Germain, 10.
Bibonne, architecte, Magasins, 12.
Bibron, aide-naturaliste au Jardin-des-Plantes, Carier, 20.
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bican (Vve) et fils, fondeurs en cuivre, place de la Corderie-du-Temple, 26.
Bicel, épicier, Marché-d'Aguesseau, 15.
Bichard, tabac et eau-de-vie, Faub.-St-Martin, 45.

Didot 1844a - pages 125-126

Un instantané des acteurs commerciaux à un instant

Des entrées **disparaissent** au cours du temps.

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33
(provisoirement).
***Bibl'que protestante (Société)*, Moulins, 16.**
***Bibron, aide-natural*, au Muséum d'Hist. nat.**
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bicant (Mme), tondeur en cuivre, cour de la
Corderie-du-Temple, 26.
Bichard (Mme), Nve-de-Luxembourg, 17.
Bichard, tabacs et eau-de-vie, Faub.-St-Mar-
tin, 45.

Didot 1841a - page 95

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33
(provisoirement).
***Bibl'que protestante (Société)*, Moulins, 16.**
***Bibron, aide-natural*, au Muséum d'Hist. nat.**
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical fab. de jouets, Montmorency, 33.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bican (Vve) et fils, tondeur en cuivre, place
de la Corderie-du-Temple, 26.
Bicel, épicier, marché d'Aguesseau, 15.
Bichard (Mme), Nve-de-Luxembourg, 17.
Bichard, tabacs et eau-de-vie, Faub.-St-Mar-
tin, 45.

Didot 1842a - page 117

Bibliothèque Ste-Geneviève, rue des Sept-Voies
et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33
(provisoirement).
Bibus, tailleur, Richelieu, 31.
Bical fab. de jouets, Montmorency, 33.
Bical et Doire, fab. de socques, Vert-Bois, 14.
Bican (Vve) et fils, fondeurs en cuivre, place
de la Corderie-du-Temple, 26.
Bicel, épicier, Marché-d'Aguesseau, 15.
Bichard, tabac et eau-de-vie, Faub.-St-Mar-
tin, 45.

Didot 1843a - page 129

Bibliothèque Ste-Geneviève, rue des Sept-Voies
et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33
(provisoirement).
Bibolet, relieur, passage Sainte-Marie-Saint-
Germain, 10.
Bibonne, architecte, Magasins, 12.
Bibron, aide-naturaliste au Jardin-des-Plan-
tes, Cuvier, 29.
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bican (Vve) et fils, fondeurs en cuivre, place
de la Corderie-du-Temple, 26.
Bicel, épicier, Marché-d'Aguesseau, 15.
Bichard, tabac et eau-de-vie, Faub.-St-Mar-
tin, 45.

Didot 1844a - pages 125-126

Un instantané des acteurs commerciaux à un instant

D'autres apparaissent.

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibl'que protestante (Société), Moulins, 16.
Bibron, aide-natural., au Muséum d'hist. nat.
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bicant (Mme), fondeur en cuivre, cour de la Corderie-du-Temple, 26.
Bichard (Mme), Nve-de-Luxembourg, 17.
Bichard, tabacs et eau-de-vie, Faub.-St-Martin, 45.

Didot 1841a - page 95

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibl'que protestante (Société), Moulins, 16.
Bibron, aide-natural., au Muséum d'hist. nat.
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bican (Vve) et fils, fondeur en cuivre, place de la Corderie-du-Temple, 26.
Bicel, épicier, marché d'Agnesseau, 15.
Bichard (Mme), Nve-de-Luxembourg, 17.
Bichard, tabacs et eau-de-vie, Faub.-St-Martin, 45.

Didot 1842a - page 117

Bibliothèque Ste-Geneviève, rue des Sept-Voies et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bical et Doire, fab. de socques, Vert-Bois, 14.
Bican (Vve) et fils, fondeurs en cuivre, place de la Corderie-du-Temple, 26.
Bicel, épicier, Marché-d'Agnesseau, 15.
Bichard, tabac et eau-de-vie, Faub.-St-Martin, 45.

Didot 1843a - page 129

Bibliothèque Ste-Geneviève, rue des Sept-Voies et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33 (provisoirement).
Bibolet, relieur, passage Sainte-Marie-Saint-Germain, 10.
Bibonne, architecte, Magasins, 12.
Bibron, aide-naturaliste au Jardin-des-Plantes, Cuvier, 29.
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bican (Vve) et fils, fondeurs en cuivre, place de la Corderie-du-Temple, 26.
Bicel, épicier, Marché-d'Agnesseau, 15.
Bichard, tabac et eau-de-vie, Faub.-St-Martin, 45.

Didot 1844a - pages 125-126

Un instantané des acteurs commerciaux à un instant

Et certaines changent.

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33
(provisoirement).
Bibl'que protestante (Société), Moulins, 16.
Bibron, aide-natural., au Muséum d'hist. nat.
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bicant (Mme), fondeur en cuivre, cour de la
Corderie-du-Temple, 26.
Bichard (Mme), Nve-de-Luxembourg, 17.
Bichard, tabacs et eau-de-vie, Faub.-St-Martin, 45.

Didot 1841a - page 95

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33
(provisoirement).
Bibl'que protestante (Société), Moulins, 16.
Bibron, aide-natural., au Muséum d'hist. nat.
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bican (Vve) et fils, fondeur en cuivre, place
de la Corderie-du-Temple, 26.
Bicel, épicier, marché d'Aguesseau, 15.
Bichard (Mme), Nve-de-Luxembourg, 17.
Bichard, tabacs et eau-de-vie, Faub.-St-Martin, 45.

Didot 1842a - page 117

Bibliothèque Ste-Geneviève, rue des Sept-Voies
et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33
(provisoirement).
Bibus, tailleur, Richelieu, 31.
Bical, fab. de jouets, Montmorency, 33.
Bical et Doire, fab. de socques, Vert-Bois, 14.
Bican (Vve) et fils, fondeurs en cuivre, place
de la Corderie-du-Temple, 26.
Bicel, épicier, Marche-d'Aguesseau, 15.
Bichard, tabac et eau-de-vie, Faub.-St-Martin, 45.

Didot 1843a - page 129

Bibliothèque Ste-Geneviève, rue des Sept-Voies
et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33
(provisoirement).
Bibolet, relieur, passage Sainte-Marie-Saint-Germain, 10.
Bibonne, architecte, Magasins, 12.
Bibron, aide-naturaliste au Jardin-des-Plantes, Cuvier, 29.
Bibus, tailleur, Richelieu, 31.
Bicel, fab. de jouets, Montmorency, 33.
Bican (Vve) et fils, fondeurs en cuivre, place
de la Corderie-du-Temple, 26.
Bicel, épicier, Marché-d'Aguesseau, 15.
Bichard, tabac et eau-de-vie, Faub.-St-Martin, 45.

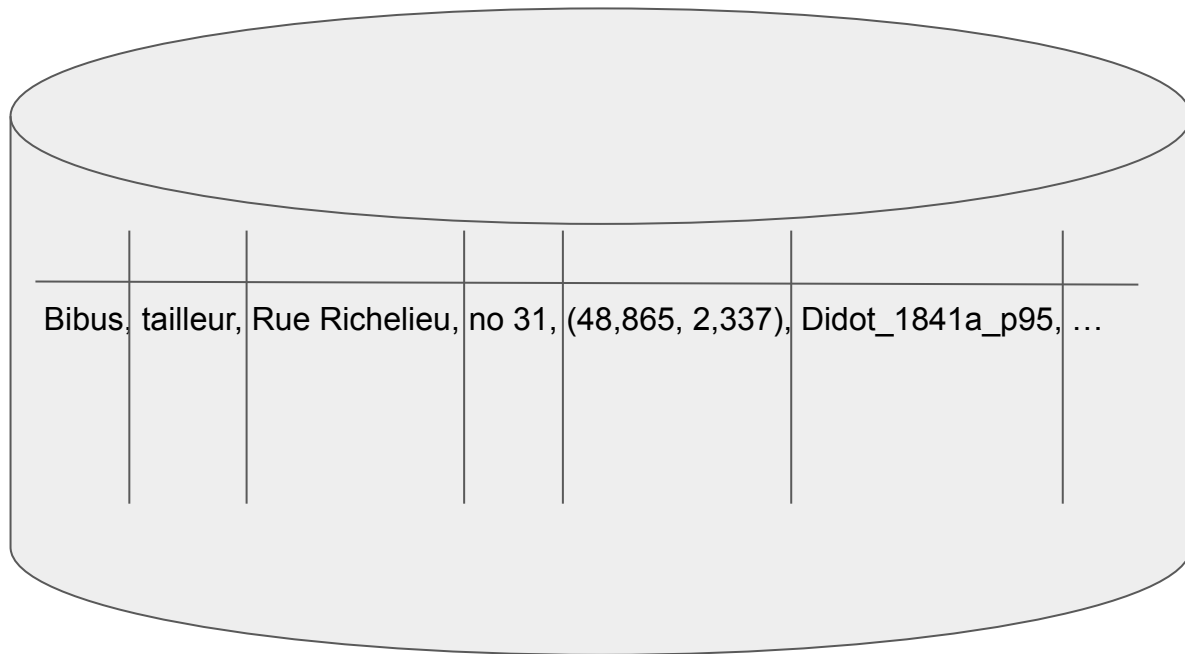
Didot 1844a - pages 125-126

Notre objectif d'extraction de données

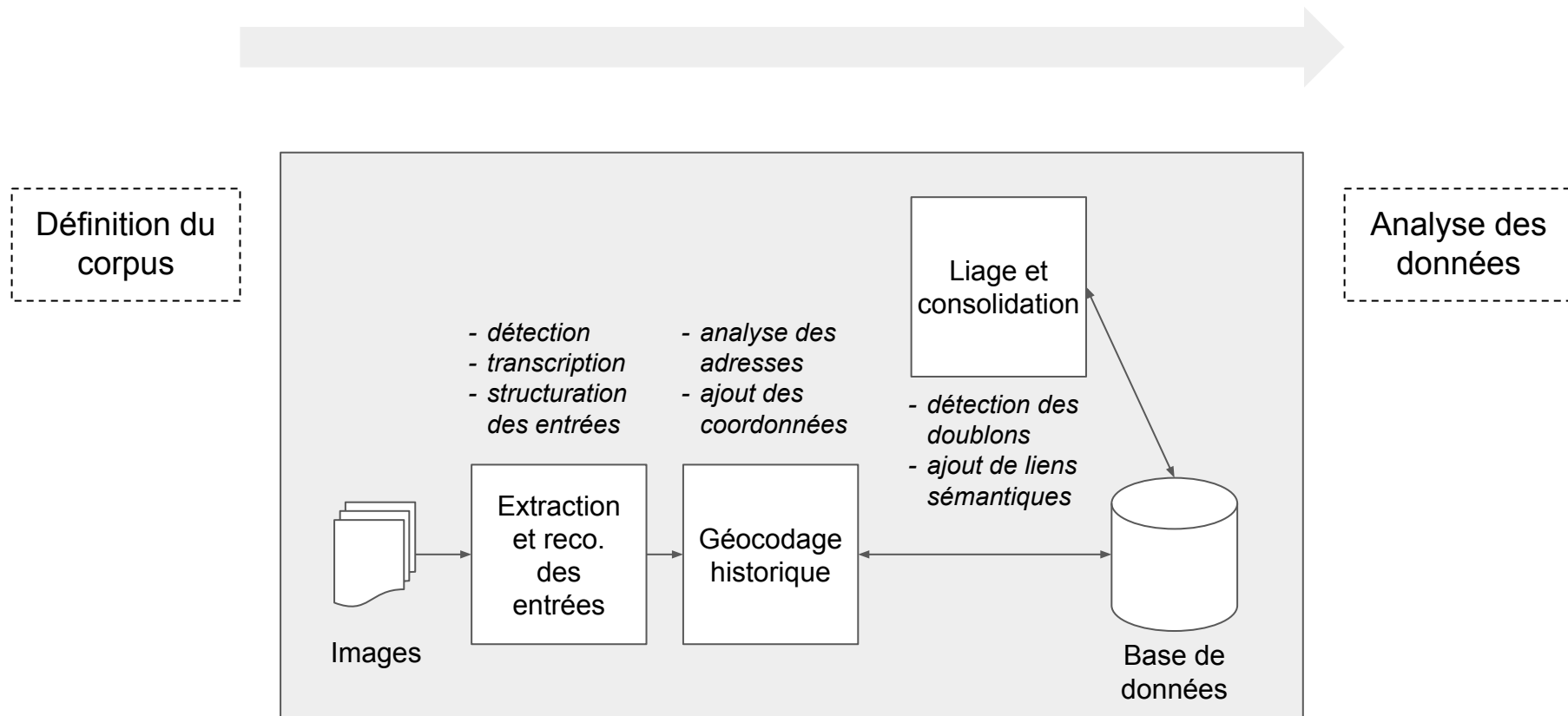
Extraire les entrées des annuaires pour les géo-référencer, les dédoublonner, les corriger, les relier...

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33
(provisoirement).
Bibl'que protestante (Société), Moulins, 16.
Bibron, aide-natural. , au Muséum d'hist. nat.
Bibus, tailleur, Roule, 21.
Bibus, tailleur, Richelieu, 31.
Bical et Dorre, fab. de socques, Vertbois, 14.
Bicant (Mme), fondateur en cuivre, cour de la
Corderie-du-Temple, 26.
Richard (Mme), Nve-de-Luxembourg, 17.
Richard, tabacs et eau-de-vie, Faub.-St-Martin, 45.

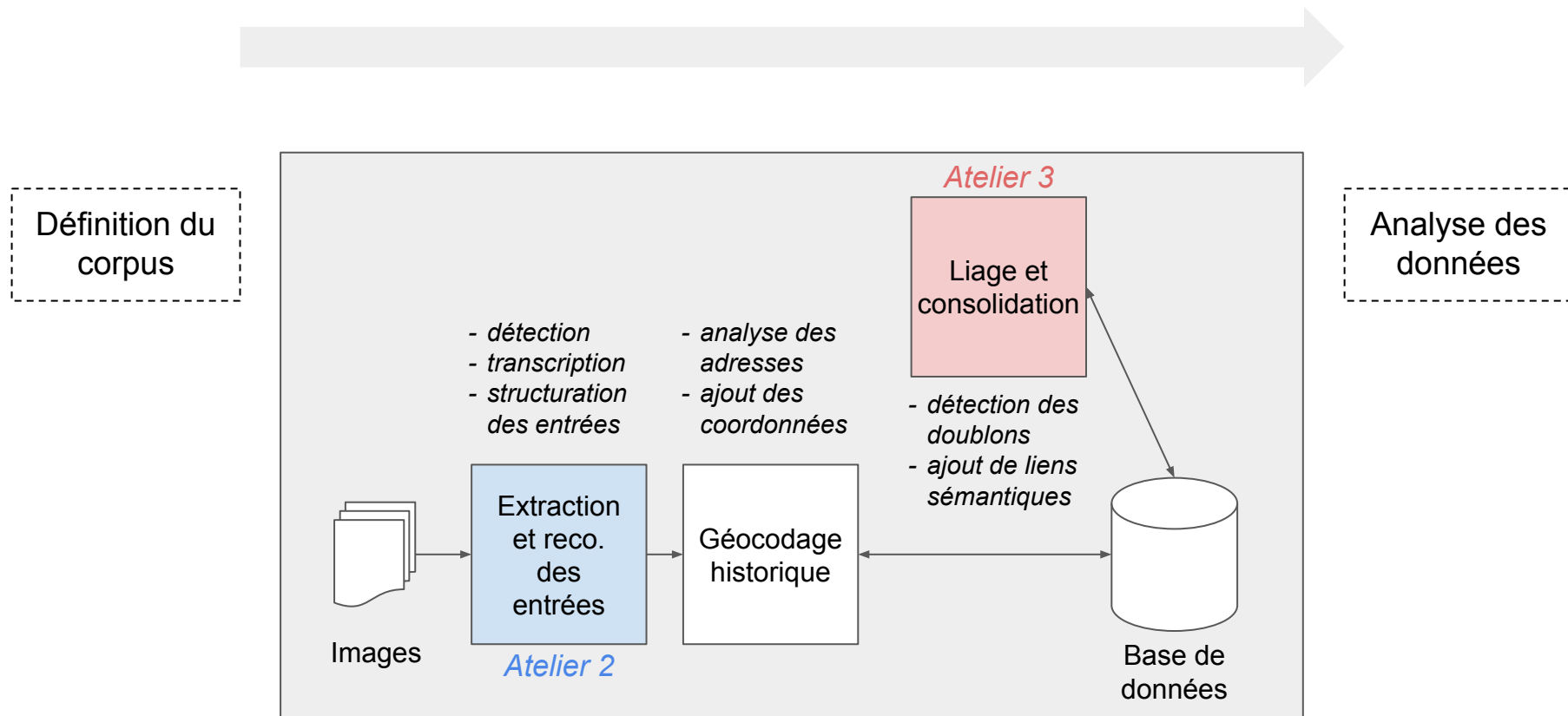
Didot 1841a - page 95



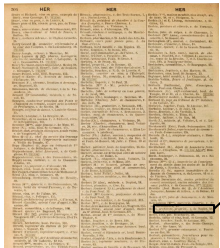
Vue d'ensemble du processus



Vue d'ensemble du processus



Extraction et reconnaissance des entrées



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



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

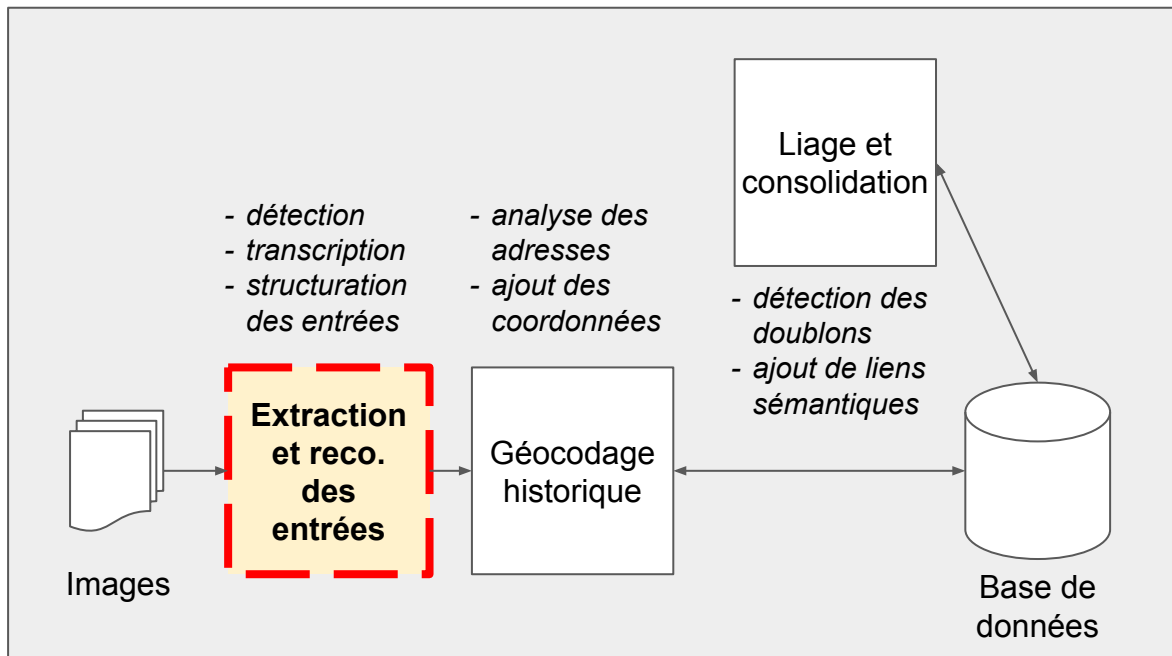


PERSONNE	ACTIVITÉ	RUE	NUMÉRO	et	RUE	NUMÉRO
Lemonnyer,	plomberie,	r. de Bondy,	86,	et	r. Bouchardon,	1.

Vue d'ensemble du processus



Définition du corpus



Analyse des données

Chaîne de traitement pour l'extract. et la reco. des entrées

Pré-traitement image
(correction géométrique,
réduction de bruit...)

Segmentation du canvas
Classification des blocs

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

Transcription

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

Extraction des
informations clés

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



Chaîne de traitement pour l'extract. et la reco. des entrées

Pré-traitement image
(correction géométrique,
réduction de bruit...)

Segmentation du canvas
Classification des blocs

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

Transcription

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

Extraction des
informations clés

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



Amélioration de la qualité d'image

□ Réduire la variance des entrées pour la suite des traitements

- ✓ Amélioration de contraste / suppression de fond
- ✓ Rotation / Shear



Une contribution notable publiée :
détection de lignes rapide et fiable

P. Bernet, J. Chazalon, E. Carlinet, A. Bourquelot et E. Puybareau, "Linear Object Detection in Document Images using Multiple Object Tracking", in Proc. ICDAR 2023

<https://github.com/EPITAResearchLab/bernet.23.icdar>

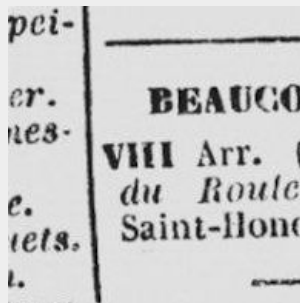
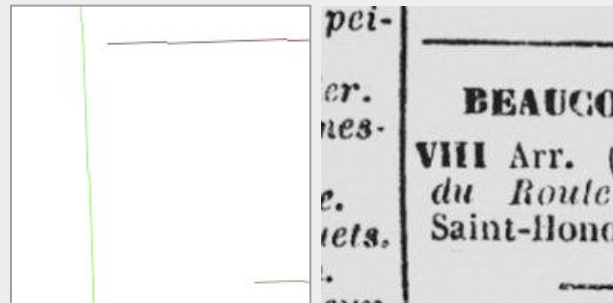


Image originale



Lignes détectées



Image redressée

Chaîne de traitement pour l'extract. et la reco. des entrées

Pré-traitement image
(correction géométrique,
réduction de bruit...)

Segmentation du canvas
Classification des blocs

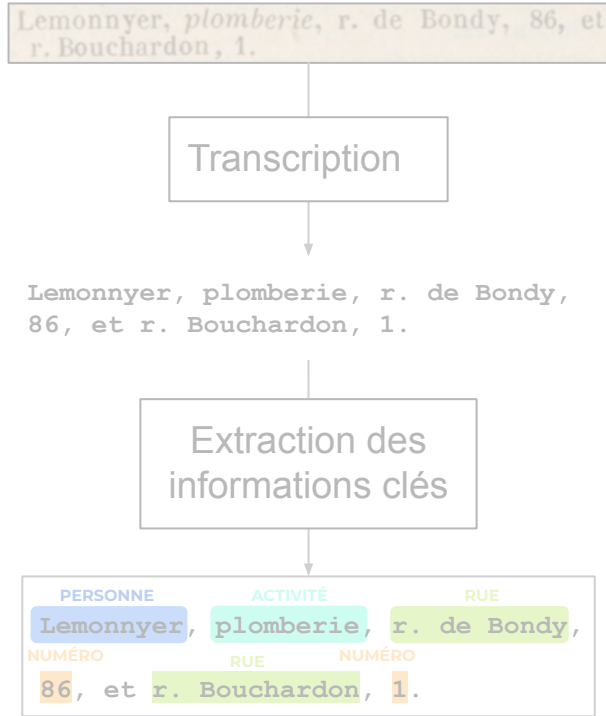
Lemonnyer, plomberie, r. de Bondy, 86, et
r. Bouchardon, 1.

Transcription

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

Extraction des
informations clés

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



Extraction de canvas / Segmentation

Méthode *ad hoc* qui tire profit de la mise en page régulière

1. Séparation en **blocs** (XY cuts, smearing) et Classification des régions (entête, titre...)
2. Détection des **lignes** (watershed)
3. Regroupement des lignes en **entrées** (HMM)

Avantages

- Très rapide (fraction de seconde par image)
- Aussi efficace que les approches modernes type LayoutParser sur ces données
- *Extraction d'entrées intégrée*

Limites

- *Extraction d'entrées intégrée*
- Pas d'exploitation de l'information textuelle
- Pas de gestion du multi-pages
- Sensible au bruit, d'où le nettoyage en amont

FORMES A SUCRE.	FOURCHES ET PELLÉS.	1410
BROMER ferret, rue des Formes 35.	BURGEL (C) G. Lagomoni et C. Avio, avenue de Vieux, 30, cochenes, mousettes pour crautes. EN GROS , Monton d'Yvon, J. Marouin & J. Petit Tronostretrassart, 13 et 15 Lombard, rue de l'Éclair, 124, Berlin, Huppelstraße.	BERGER & Cie, appareil d'écriture, fabrique de crautes, pap. National, 2, r. Pa roust, Seine à Paris, r. de l'Éclair, 124.
BUSSY cravaches, toile et rubanets blancs et blancs. Toiles et mousses, blancs, 1511, rue de la Tulipe, Valenciennes.	BERVARD F. F. J. J. J. J. J. J. J. J. J. dionne, (crautes, et cravaches en gros, rue St-Martin, 111 et 112.	FOURBISSEURS. (Voir aussi Foibles.)
CHEVALER cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.	BELLEVILLE (C) G. Lagomoni et C. Avio, avenue de Vieux, 30, cochenes, mousettes pour crautes. EN GROS , Monton d'Yvon, J. Marouin & J. Petit Tronostretrassart, 13 et 15 Lombard, rue de l'Éclair, 124, Berlin, Huppelstraße.	BUCKS (E), succe. de P. Dele court, et Bucks, Hézard, 7.
CHEVALER cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.	BLONDE (C) G. Lagomoni et C. Avio, avenue de Vieux, 30, cochenes, mousettes pour crautes. EN GROS , Monton d'Yvon, J. Marouin & J. Petit Tronostretrassart, 13 et 15 Lombard, rue de l'Éclair, 124, Berlin, Huppelstraße.	BUSSY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.
CHIGNY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.	BLOCH (C) G. Lagomoni et C. Avio, avenue de Vieux, 30, cochenes, mousettes pour crautes. EN GROS , Monton d'Yvon, J. Marouin & J. Petit Tronostretrassart, 13 et 15 Lombard, rue de l'Éclair, 124, Berlin, Huppelstraße.	BUSSY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.
CHIGNY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.	BOISSE (C) G. Lagomoni et C. Avio, avenue de Vieux, 30, cochenes, mousettes pour crautes. EN GROS , Monton d'Yvon, J. Marouin & J. Petit Tronostretrassart, 13 et 15 Lombard, rue de l'Éclair, 124, Berlin, Huppelstraße.	BUSSY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.
CHIGNY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.	BOISSE (C) G. Lagomoni et C. Avio, avenue de Vieux, 30, cochenes, mousettes pour crautes. EN GROS , Monton d'Yvon, J. Marouin & J. Petit Tronostretrassart, 13 et 15 Lombard, rue de l'Éclair, 124, Berlin, Huppelstraße.	BUSSY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.
CHIGNY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.	BOISSE (C) G. Lagomoni et C. Avio, avenue de Vieux, 30, cochenes, mousettes pour crautes. EN GROS , Monton d'Yvon, J. Marouin & J. Petit Tronostretrassart, 13 et 15 Lombard, rue de l'Éclair, 124, Berlin, Huppelstraße.	BUSSY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.
CHIGNY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.	BOISSE (C) G. Lagomoni et C. Avio, avenue de Vieux, 30, cochenes, mousettes pour crautes. EN GROS , Monton d'Yvon, J. Marouin & J. Petit Tronostretrassart, 13 et 15 Lombard, rue de l'Éclair, 124, Berlin, Huppelstraße.	BUSSY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.
CHIGNY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.	BOISSE (C) G. Lagomoni et C. Avio, avenue de Vieux, 30, cochenes, mousettes pour crautes. EN GROS , Monton d'Yvon, J. Marouin & J. Petit Tronostretrassart, 13 et 15 Lombard, rue de l'Éclair, 124, Berlin, Huppelstraße.	BUSSY cravaches, toile et rubanets blancs, 1511, rue de la Tulipe, Valenciennes.

Chaîne de traitement pour l'extract. et la reco. des entrées

Pré-traitement image
(correction géométrique,
réduction de bruit...)

Segmentation du canvas
Classification des blocs

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

Transcription

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

Extraction des
informations clés

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



Reconnaissance du texte (OCR)

Comparaison de plusieurs systèmes OCR sur notre jeu de données “A Dataset of French Trade Directories from the 19th Century (FTD)”

DOI [10.5281/zenodo.6394464](https://doi.org/10.5281/zenodo.6394464)

Systèmes testés

- [Tesseract v4](#) (v5 pas testée)
- [PERO OCR](#) (code Github + modèle 2020 auteurs)
- [Kraken OCR](#) (modèle générique EN imprimé)

Pas de ré-entraînement.

Bottin 1820

Dufort, *bottier*, Palais-R., gal. vitrée, 215.
295

Bottin 1827

Baleste, *chef aux domaines*, S.-Georges, 17.

Bottin 1837

Cattois, *pharmac.*, Bretagne, 46.

Bottin 1854

Fontaine, *draperies*, Neuve-des-Petits-Champs, 2.

Cambon Almgene 1841

Aron Javal (L.) art. de Paris, r. des Bourdonnais, 17.

Deflandre 1828

DEVILLERS, r. Croix-des-Pet.-Champs, 25.
Cordonn.

Deflandre 1829

Huguenin, *épïc.*, r. de Valois, 8, Pal.-Royal.

Didot 1851

Viéville, *fab. de boutons*, Aumaire, 48, et place St-Nicolas-des-Champs, 2.

Reconnaissance du texte (OCR)

Résultats

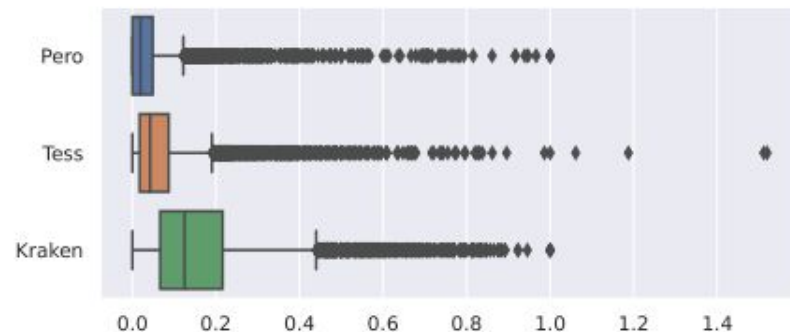
- Très bonne performance de PERO OCR
- Manque de modèles publics pour Kraken (en évolution avec [HTR-United](#))

Autres systèmes libres compétitifs testés depuis

- [Microsoft UniLM/TrOCR](#)
- [Mindee docTR](#)

⇒ OCR plutôt performant, résultats exploitables

	PERO OCR	Tesseract	Kraken
CER	3.78%	6.56%	15.72%



Chaîne de traitement pour l'extract. et la reco. des entrées

Pré-traitement image
(correction géométrique,
réduction de bruit...)

Segmentation du canvas
Classification des blocs

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

Transcription

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

Extraction des
informations clés

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



Reco. d'entités nommées (*Named Entity Recognition*)

Identifier et classer des expressions d'un ou plusieurs mots d'un texte en catégories prédéfinies : personne, lieu, organisation, date, prix...

Entités spécifiques dans le cas des annuaires.

Exemple:

PERSON
Ravrio et comp., **ACTIVITY**
fabr. de bronzes et curiosités,
r. Richelieu, **CARDINAL** 93; **FEATURE TYPE** la fabrique **LOCATION** rue Montmartre, **CARDINAL** 161.

Difficultés : Résultats OCR bruités

Gavarret ✱, prof. de physique à la faculté de
médecine, Grenelle-St-Germain, 49.


ravarret #, prof. de physique à la faculté de
médecine, Grenelle-St-Germain, 49.

Tesseract v4 output

Duffaut, chaudronnier, r. de la Sourdière,
314

Daflant, c'audronnier, v, de la RARES
Ge OO a x

Tesseract v4 output

 = erreurs

Besoin de **robustesse au bruit OCR**, voire de *post-correction*.

Difficultés : Variation de la structure des entrées

Raison soc.

Activité

Adresse

Durand jeune; pour bas, Charenton, 12 ancien. 18. *

No rue

Prévost-Guillaume, f. ta., r. N.-St.-Mart., 28.

Lefranc Méquignon et co., satins turcs, prunelle, satins et draps de soie, gros de Naples, toiles, coutils, galons, rubans, coulisses et lacets pour chaussures de dames, sommières, flanelles et molletons de soie pour fourrures, r. des Prouvaires, 32.

Planche , R. de Poitou, 9.-H. Armé.

Baronnat frères, soies teintes et écrues, fil-Denis, 257, passage du Renard ; maison à Lyon, r. Cen-

Jamain, orangiste, ⒶS. H.1831, Fos-sés-St-Marcet, 12.

Appert fils, verres et cristaux, 21-23 Jour.

Mabire, Lourcine, 124.

Besoin d'un système robuste aux **variations** et **bruit syntaxique**, via un **apprentissage** à base d'**exemples**.

Approche NER face aux données bruitées

Utilisation d'un modèle NER Transformer **CamemBERT** déjà entraîné, et spécialisation à nos données :

- Pré-entraînement sur données OCR brutes (bruitées)
- Entraînement supervisé sur données OCR bruitées alignées sur la référence

⇒ Gains significatifs en performance, F-score de détection > 94% sur test set

Publication associée :

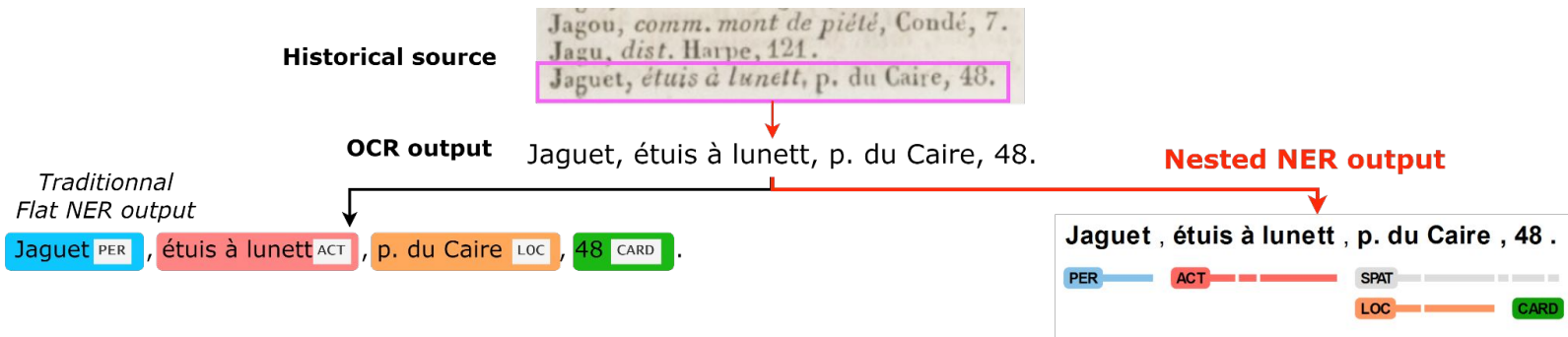
N. Abadie, E. Carlinet, J. Chazalon et B. Duménieu, "A Benchmark of Named Entity Recognition Approaches in Historical Documents Application to 19th Century French Directories", in Proc. DAS 2022, <https://github.com/soduco/paper-ner-bench-das22>

Expérience NER imbriqué (non déployée)

Publication associée :

S. Tual, N. Abadie, E. Carlinet, J. Chazalon, et B. Duménieu, "A Benchmark of Nested NER Approaches in Historical Structured Documents", in Proc. ICDAR 2023, <https://github.com/soduco/paper-nestedner-icdar23-code/>

⇒ Extraction riche possible, performance inchangée



En pratique...

[Extraction “v2”](#), sur ~110 ouvrages, début 2023

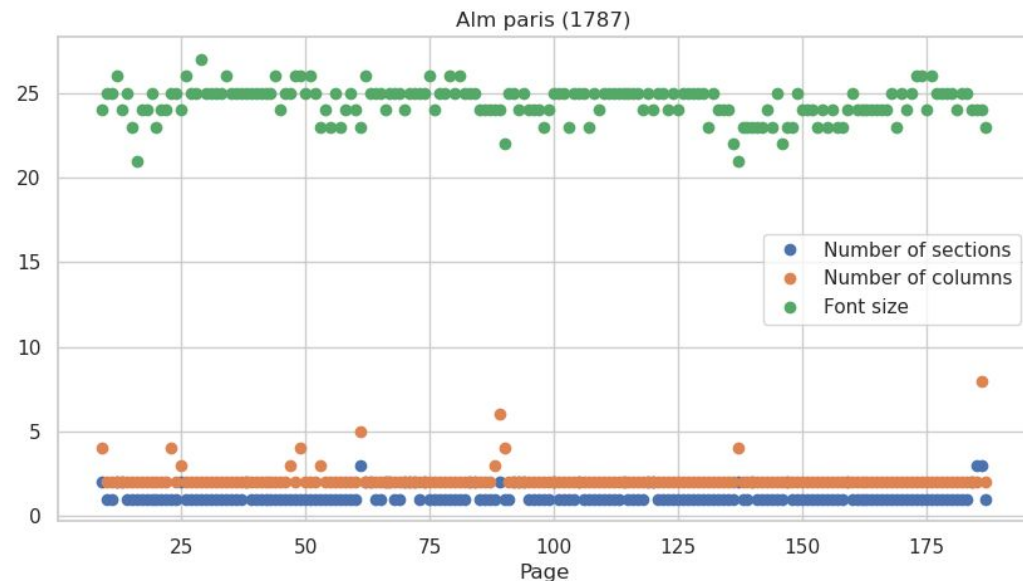
- 9 825 018 entrées extraites
- Environ 74% d’entrées exploitables, de la forme :

PERSON [**ACT**] **LOC** [**CARDINAL**]

Amélioration 1 : Validation de la cohérence inter-pages

Mécanisme en 2 passes :

1. Traiter les images séparément, et calculer certains indicateurs clés : taille des caractères, nombre de colonnes...
2. En déduire les paramètres à appliquer pour les images incohérentes ; relancer les traitements pour les pages / images concernées



Amélioration 2 : Traitement du texte en flux (en cours)

Problèmes attaqués :

1. Comment tirer profit à la fois de l'**information visuelle et textuelle** pour mieux **séparer les entrées** ?
2. Comment s'affranchir des **sauts de colonne** et de **pages** ?
3. Est-il possible de **séparer les entrées** et faire le **NER simultanément** ?

Opportunité : documents majoritairement textuels, ordre de lecture simple

Solution : un modèle de langage mixte **visuel** et **textuel**

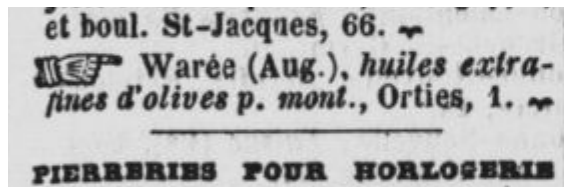


Image originale

```
<page_break><space> et boul. St-Jacques, 66. <space>  
<line_break> <space>☞ Warée (Aug.), huiles extra-  
<line_break> <space> fines d'olives p. mont., Orties, 1. <space>  
<line_break> <space> — <space>  
<line_break> <space> PIERRERIES POUR HORLOGERIE <space>
```

Flux d'entrée enrichi pour le système "NER+sép. entrées"

Amélioration 2 : Traitement du texte en flux (en cours)

Capacité à extraire des entrées à cheval sur plusieurs pages attirante.

Mais une tendance à fusionner les entrées.

⇒ Des résultats intéressants mais à améliorer.

- À combiner avec une approche exploitant la position des éléments type LayoutLM.
- Travaux en cours.

Au final...

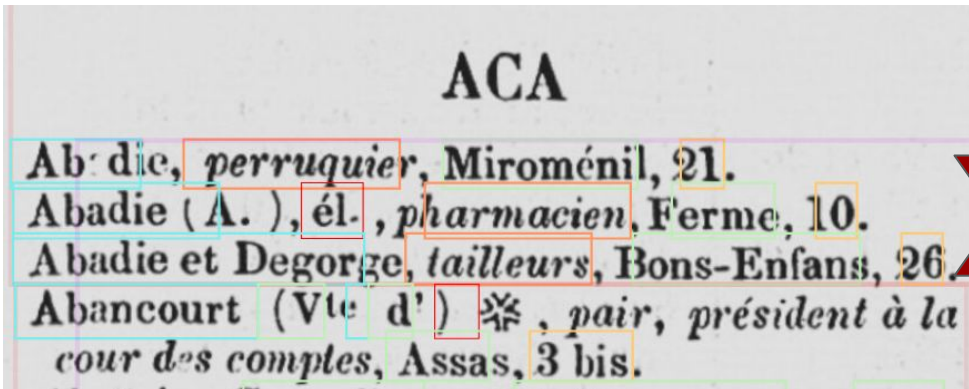
Extraction “v4”, avec séparation simple des entrées fusionnées et réassociation “nom de voie, numéro”, sur 357 listes réparties entre 144 ouvrages et 192 PDFs :

- 22 743 928 entrées
- 23 728 378 entités PER
- 2 106 808 entités TITRE
- 14 663 426 entités ACT
- 27 618 637 addresses (combinaisons de CARDINAL + LOC, ou LOC seules)
dont 96.3% ont été géocodées

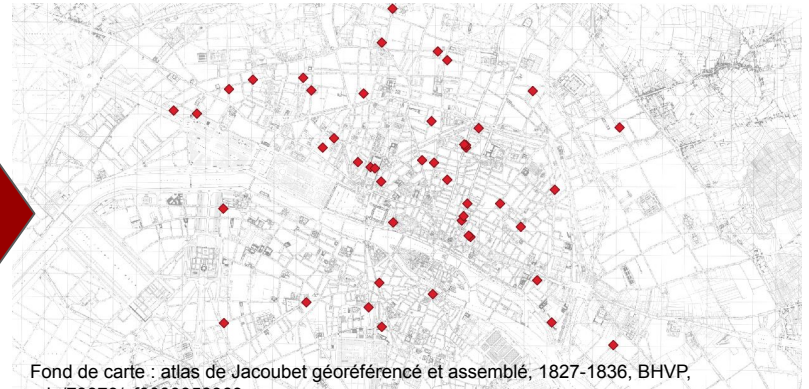
Soit environ 27 millions de points à placer sur une carte !

Géocodeur historique

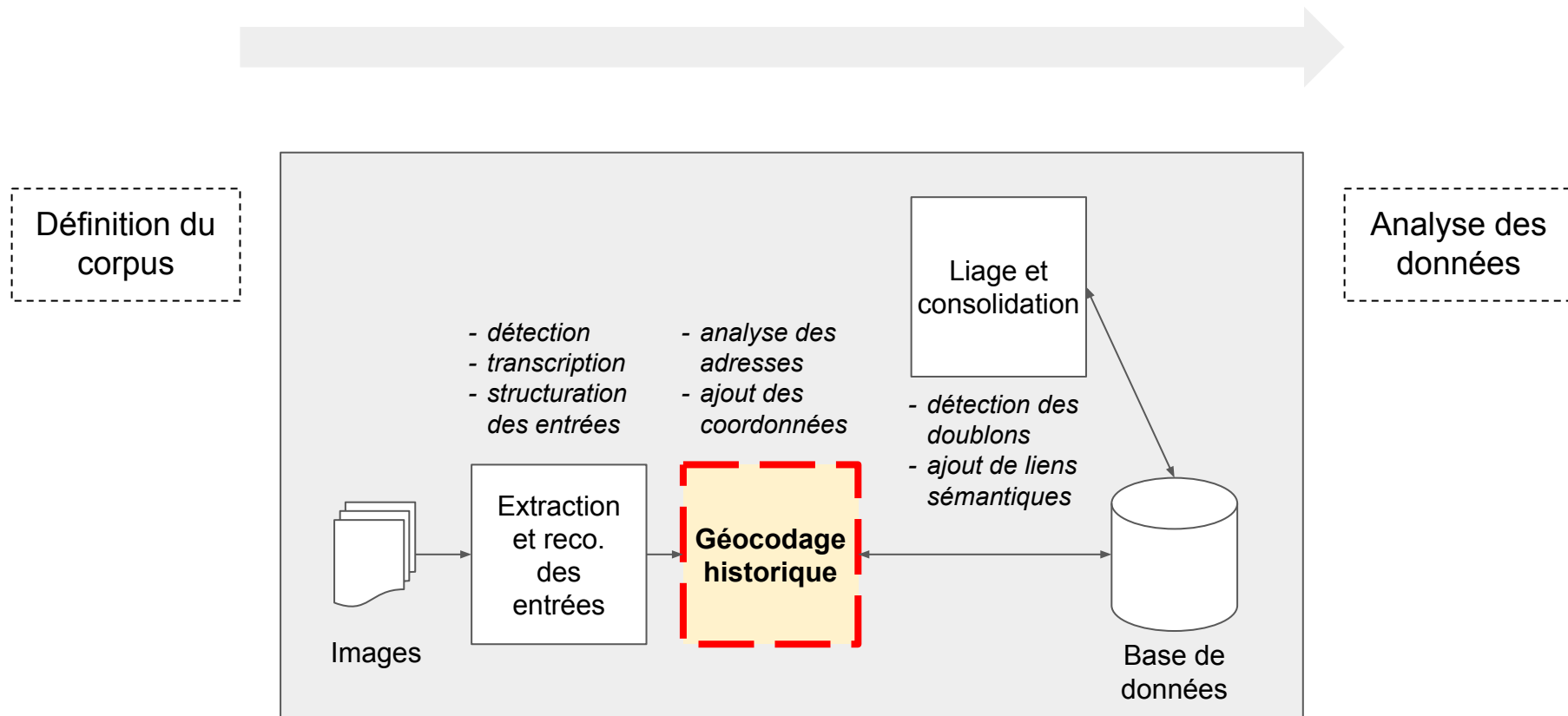
(Lamy_1840) Annuaire général du commerce, judiciaire et administratif de France et des principales villes du monde, v. 128.



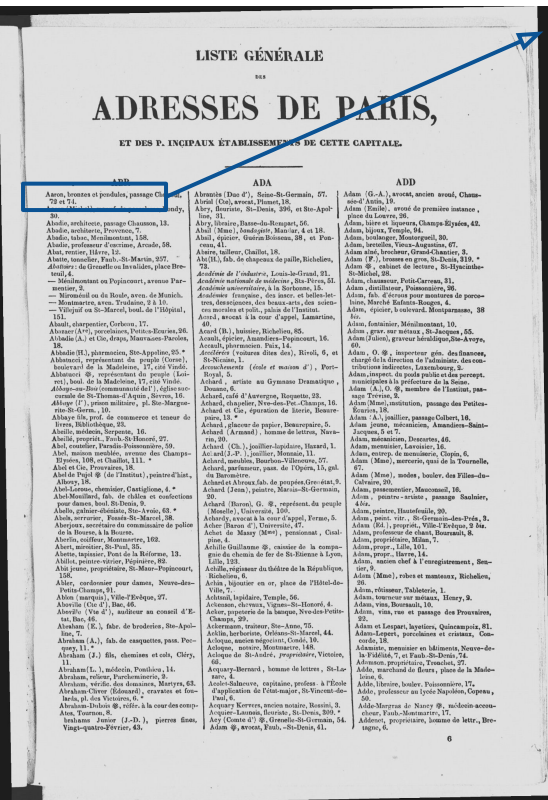
https://api.geohistoricaldata.org/directories/export_geojson?source.pdf_id=eq.Lamy_1840&source.pdf_view=eq.0128&



Vue d'ensemble du processus



Chaîne de traitement



#1 OCR

> Aaron, bronzes et pendules, passage Choiseul, 72 et 74.

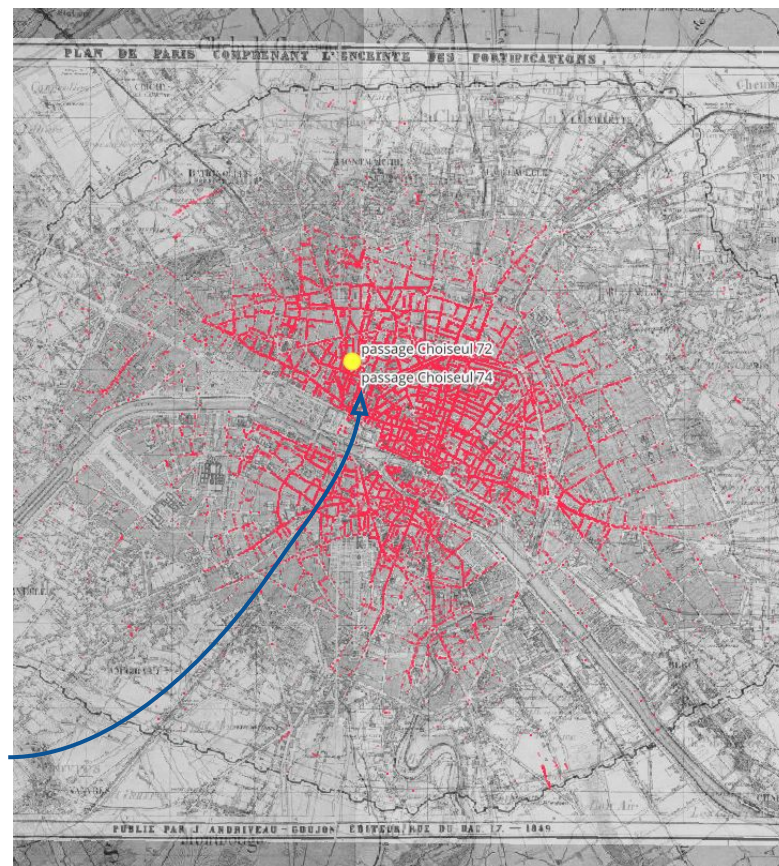
#2 NER

> Aaron, bronzes et pendules, passage Choiseul, 72 et 74.

#3 Restructuration

Address #1 : passage Choiseul 72

Address #2 : passage Choiseul 74



Géocodeur historique: un moteur de recherche exploitant les données des atlas de Paris

adresse
"passage Choiseul, 72"

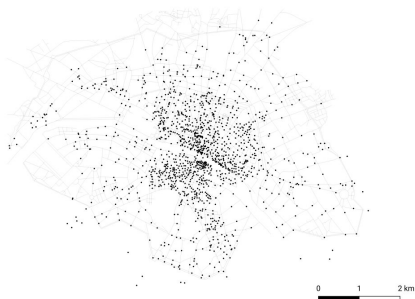
année visée (scalaire ou intervalle)
[1850,1852]

Recherche multicritères dans les données vectorielles des plans :

- matching flou du nom + numéro de rue
- distance temporelle



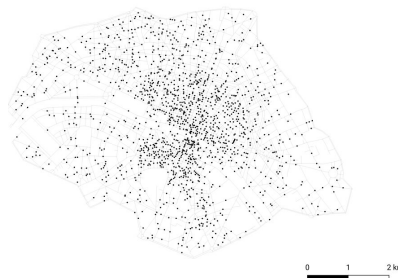
Adaptation géohistorique du géocodeur libre et *open source* PELIAS



Atlas de Verniquet, 1784-1795 ⌚
Filaire des rues



Atlas de Jacoubet, 1827-1836 ⌚
Filaire des rues & points adresse



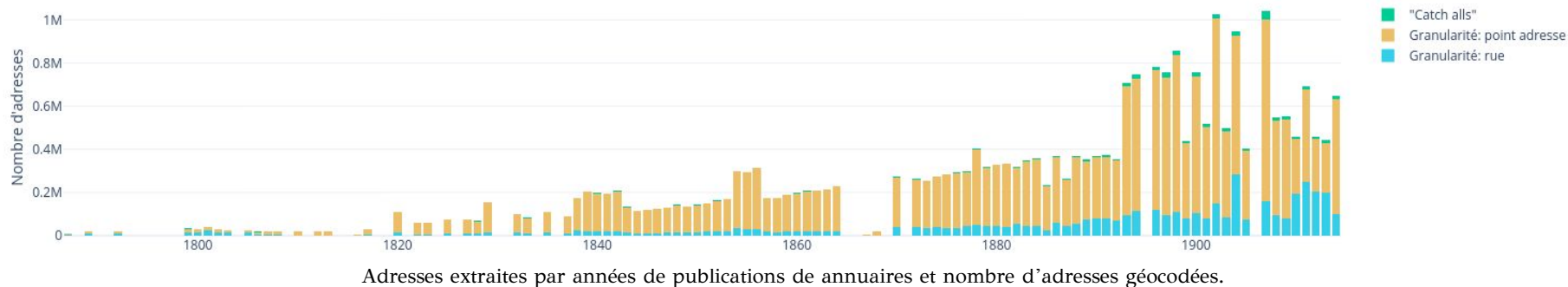
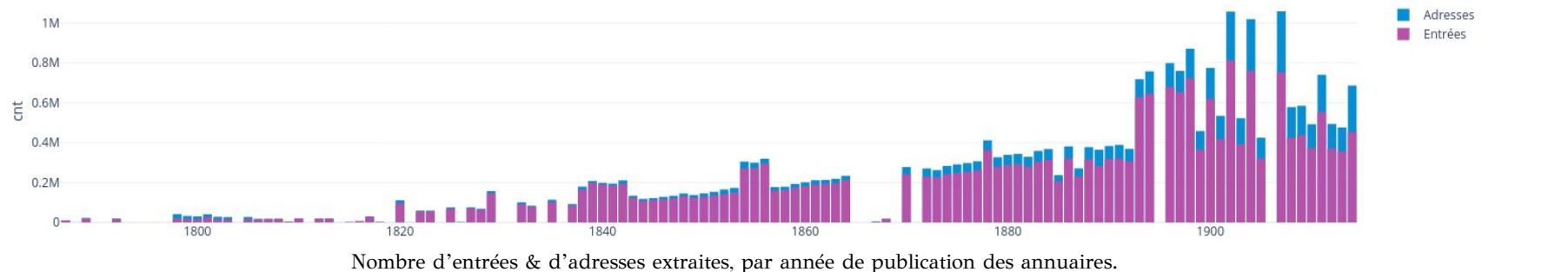
Plan Andriveau-Goujon, 1849 ⌚
Filaire des rues



Atlas municipal des 20 arrds., 1888 ⌚
Filaire des rues & points adresse

Géocodage des annuaires : aperçu des résultats

V2, octobre 2023 : 96% d'entrées géocodées (V2 septembre 2022 : 66%)

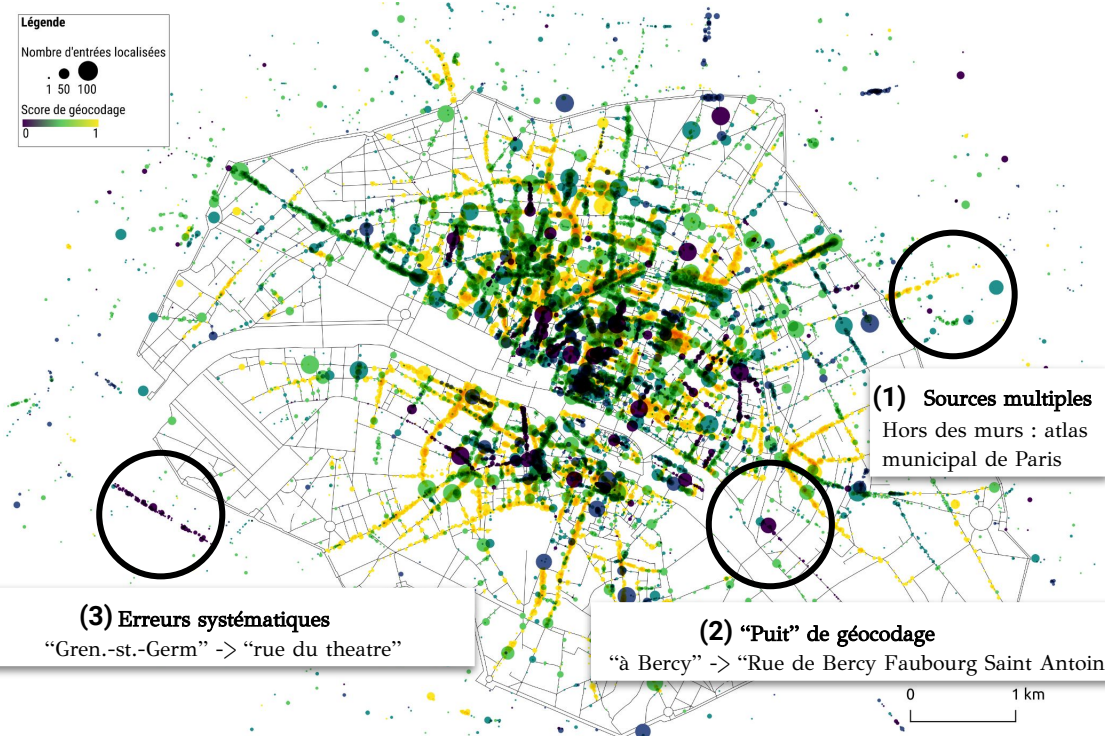
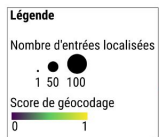


Résultats du géocodage des annuaires



Entrées géocodées pour l'année 1838 : Almanach du commerce de Paris[...] (Bottin) et Annuaire général du commerce, de l'industrie et de l'agriculture de France [...] (Henrichs).

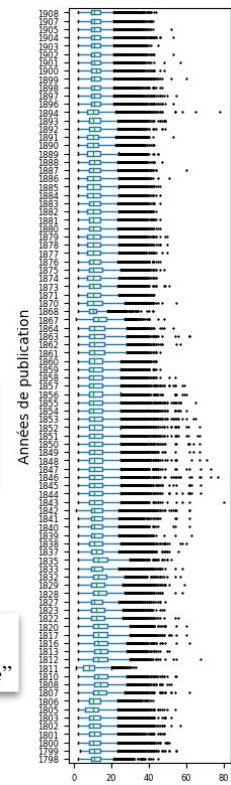
Longueurs des énoncés d'adresses en nombre de caractères et exemples typiques



(1) Sources multiples
Hors des murs : atlas municipal de Paris

(2) "Puit" de géocodage
"à Bercy" -> "Rue de Bercy Faubourg Saint Antoine"

(3) Erreurs systématiques
"Gren.-st.-Germ" -> "rue du theatre"



Nombre de caractères

r. d. des, 123, che 38,
Bar 63, NAS, r. d, bis,
Tri, imp, BaC 104, Dou, St-, omb,
Bac 142., r.,
e N T, L Ver

Mauvaises-Paroles 15

quai aux Fleurs, au coin de la rue
de la Cité

avenue de Neuilly, sur la pelouse
près la rue du chemin de Versailles



Liage et consolidation

même personne

succession

déménagement

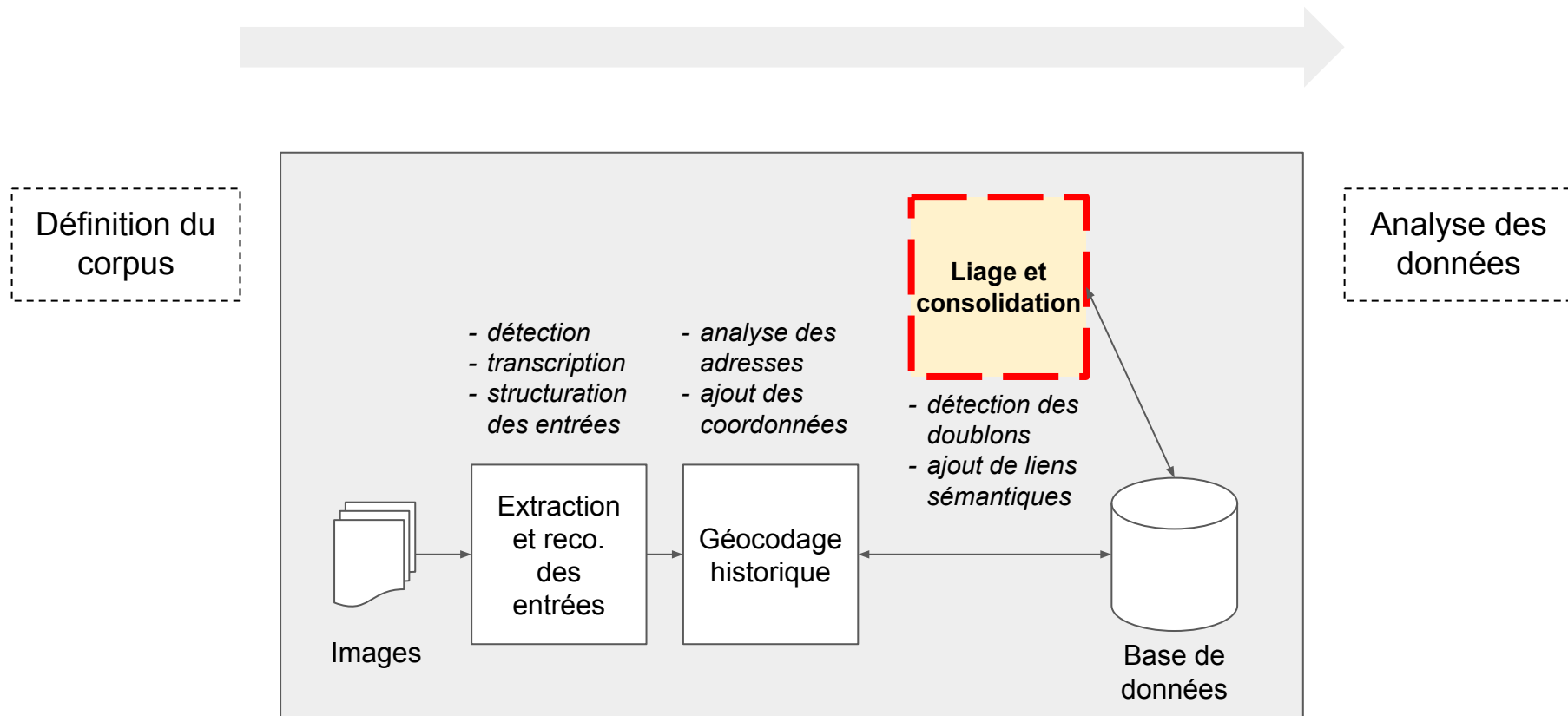
Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33
 (provisoirement).
Bibl'que protestante (Société), Moulins, 16.
 Bibron, aide-natural., au Muséum d'hist. nat.
 Bibus, tailleur, Roule, 21.
 Bibus, tailleur, Richelieu, 31.
 Bical et Dorre, fab. de socques, Vertbois, 14.
 Bican (Mme), fondeur en cuivre, cour de la
 Corderie-du-Temple, 26.
 Bichard (Mme), Nve-de-Luxembourg, 17.
 Bichard, tabacs et eau-de-vie, Faub.-St-Martin,
 45.

Bibliothèque Ste-Geneviève, Clotilde, 1.
Bibliothèque de la Ville, quai d'Austerlitz, 33
 (provisoirement).
Bibl'que protestante (Société), Moulins, 16.
 Bibron, aide-natural., au Muséum d'hist. nat.
 Bibus, tailleur, Roule, 21.
 Bibus, tailleur, Richelieu, 31.
 Bical, fab. de jouets, Montmorency, 33.
 Bical et Dorre, fab. de socques, Vertbois, 14.
 Bican (Vve) et fils, fondeur en cuivre, place
 de la Corderie-du-Temple, 26.
 Bichel, épicier, marché d'Aguesseau, 15.
 Bichard (Mme), Nve-de-Luxembourg, 17.
 Bichard, tabacs et eau-de-vie, Faub.-St-Martin,
 45.

Bibliothèque Ste-Geneviève, rue des Sept-Voies
 et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33
 (provisoirement).
 Bibus, tailleur, Richelieu, 31.
 Bical, tab. de jouets, Montmorency, 33.
 Bical et Doire, fab. de socques, Vert-Bois, 14.
 Bican (Vve) et fils, fondeurs en cuivre, place
 de la Corderie-du-Temple, 26.
 Bichel, épicier, Marche-d'Aguesseau, 15.
 Bichard, tabac et eau-de-vie, Faub.-St-Martin,
 45.

Bibliothèque Ste-Geneviève, rue des Sept-Voies
 et place du Panthéon.
Bibliothèque de la Ville, quai d'Austerlitz, 33
 (provisoirement).
 Bibolet, relieur, passage Sainte-Marie-Saint-
 Germain, 10.
 Bibonne, architecte, Magasins, 12.
 Bibron, aide-naturaliste au Jardin-des-Plan-
 tes, Carier, 20.
 Bibus, tailleur, Richelieu, 31.
 Bichel, fab. de jouets, Montmorency, 33.
 Bican (Vve) et fils, fondeurs en cuivre, place
 de la Corderie-du-Temple, 26.
 Bichel, épicier, Marché-d'Aguesseau, 15.
 Bichard, tabac et eau-de-vie, Faub.-St-Martin,
 45.

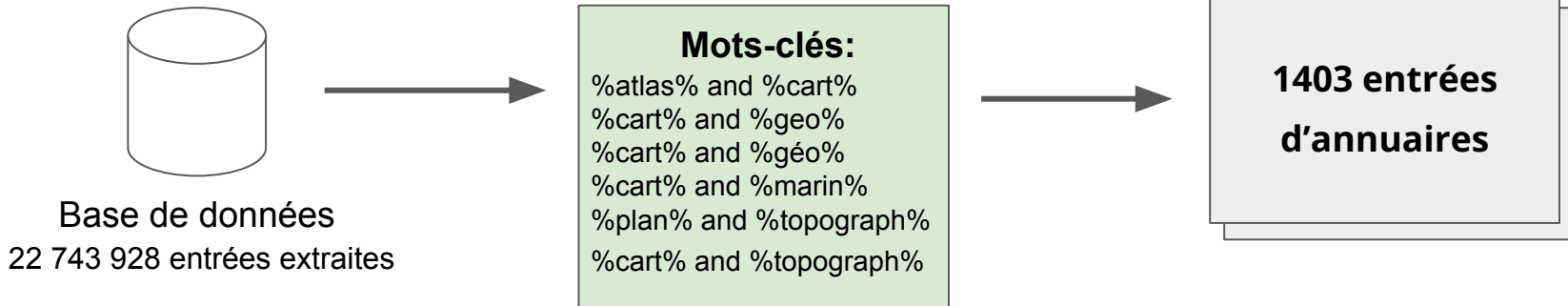
Vue d'ensemble du processus



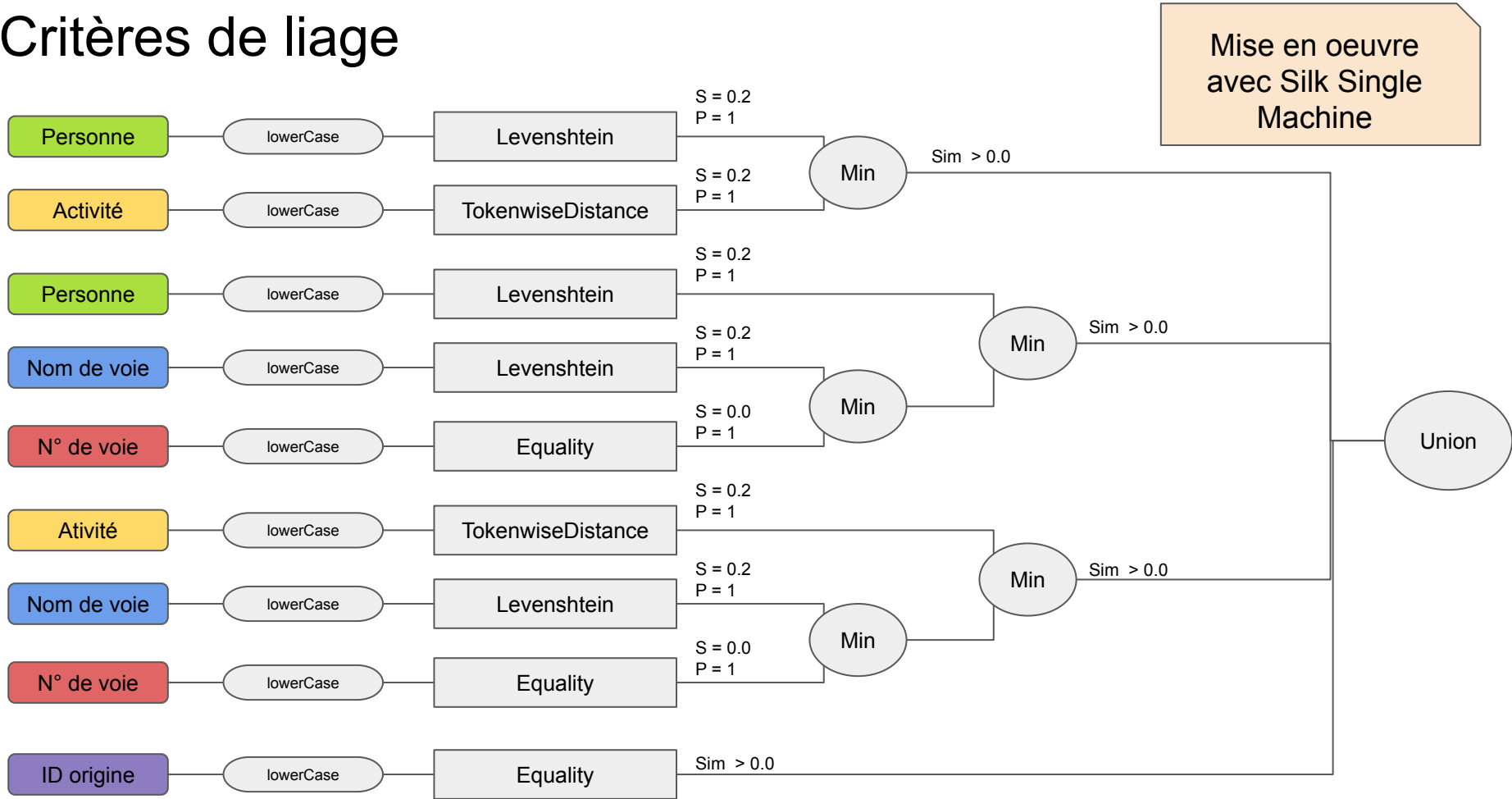
Création de graphes de connaissances géohistoriques professionnels

1. *Sélection d'un sous-ensemble d'entrées d'annuaires dans la base de données à l'aide d'une liste de mots-clés relatifs à un type d'activité.*
2. *Création de ressources RDF à partir des enregistrements de la base de données.*
3. *Utilisation d'une méthode de liage numérique pour comparer les entrées d'annuaires:*
 - *Noms similaires, activités similaires, adresses similaires → continuation*
 - *Noms similaires, activités similaires, adresses différentes → déménagement*
 - *Noms différents, activités similaires, adresses similaires → succession*
4. *Visualisation du graphe de connaissances géohistorique et évaluation qualitative des liens créés.*

Ex. Sélection des entrées sur les graveurs et marchands de cartes



Critères de liage



Résultats du liage des entrées

Exemple des graveurs et marchands de cartes

Méthode numérique

20 362 liens

Paramétrage complexe :
identifier les seuils de tolérance
pertinents

+

Raisonnement

16 046 liens

Propagation des liens
owl:sameAs par transitivité

Comparaison adaptée aux chaînes de
caractères résultant de l'OCR

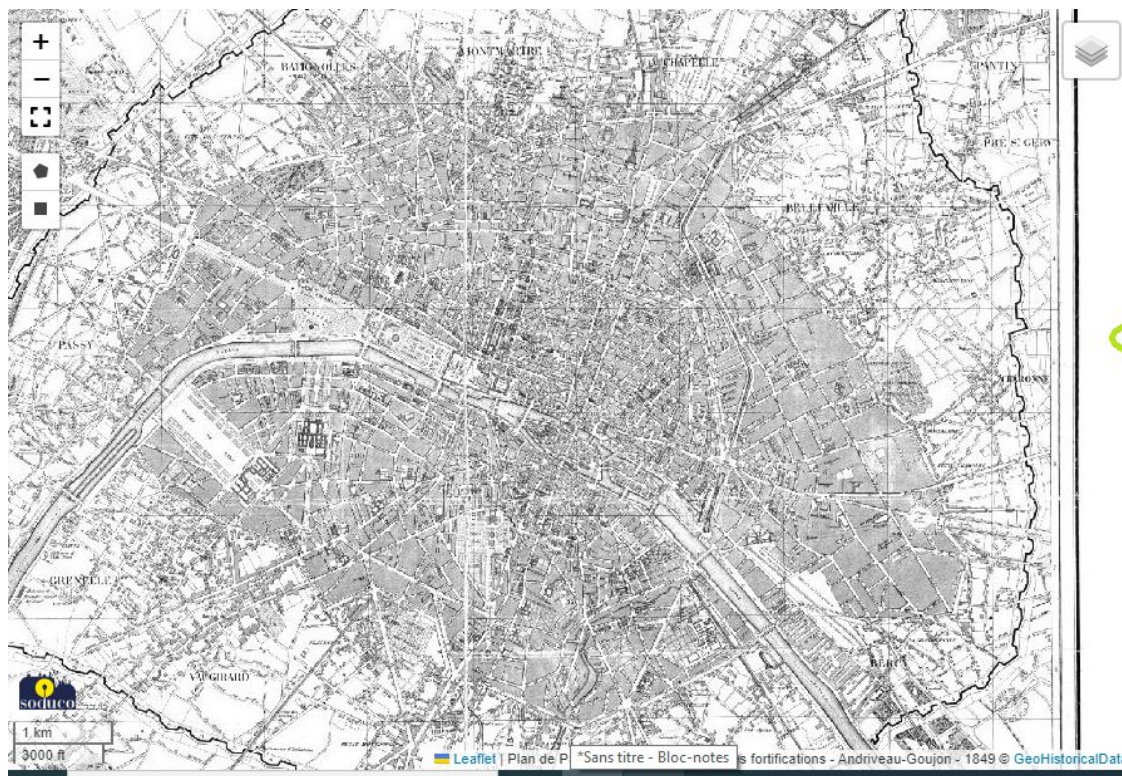
BILAN



36 408 liens d'équivalence distincts entre les ressources du graphe

Graphe géohistorique : <https://dir.geohistoricaldata.org/sparql>

Visualisation et interprétation des résultats



Visualisation des graphes géohistoriques construits à partir des entrées d'annuaires du commerce de Paris (XIXème siècle)

[Aide ?](#)

Dataset

Graveurs et marchands de cartes et plans ▾

Statistiques du dataset

Filtres

Propriétés

Raison sociale

Ex : nadar

Description

Ex : photo

Adresse

Ex : rivoli

Période

Visualisation et interprétation des résultats

Ex : nadar
Description
Ex : photo
Adresse
Ex : rivoli

Période

1840 1890

1790 1821 1853 1884 1915

? Le filtre temporel permet de faire varier l'affichage des points préalablement chargés sur la carte sans lancer une nouvelle recherche. Données chargées pour la période **1840-1890**.

Localisation

Dessinez l'emprise de votre zone de recherche avec l'outil de dessin disponible sur la carte.

Lancer la recherche

1 km
3000 ft

Leaflet | Plan de Paris contenant l'enceinte des fortifications - Andriveau-Goujon - 1849 © GeoHistoricalData

Statistiques sur le graphe géohistorique

Statistiques du jeu de données "Graveurs et marchands de cartes et plans"

Nombre d'entrées d'annuaires

1403

Nombre de ressources RDF

1732

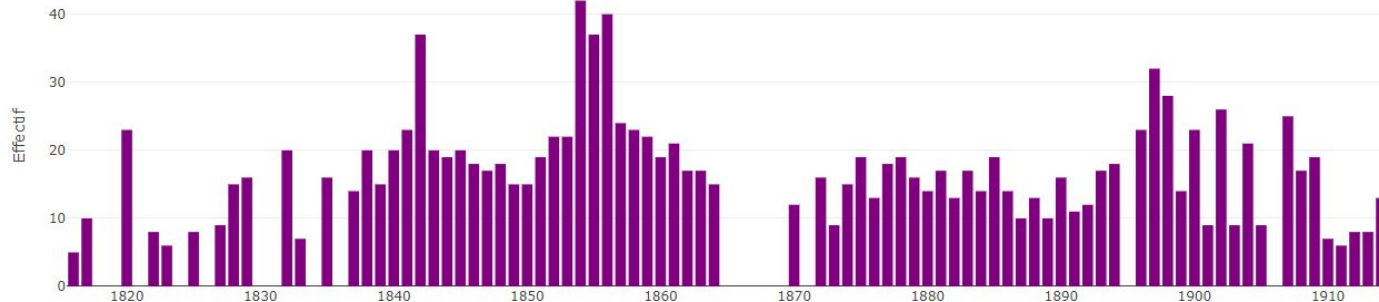
Nombre de triplets RDF

80969

Nombre de liens sameAs entre des ressources différentes

36408

Nombre d'entrées d'annuaires extraites par année



Suivi individuel des commerces



The image shows a historical map of Paris with a pop-up window for a specific location. The window contains the following information:

- Vuillemin**
- Adresse (annuaire) : 5 St-Thomas-d'Enfer
- Adresse (géocodeur) : 5 Rue Saint Thomas, Paris (Source: atlas_jacoubet_1836)
- Activité : cartes géographiques
- Année de publication : 1861
- Annuaire : DidotBottin_1861
- Identifiant de l'entrée : 1d8338bc-8614-54cc-afd5-0d412ee29a35

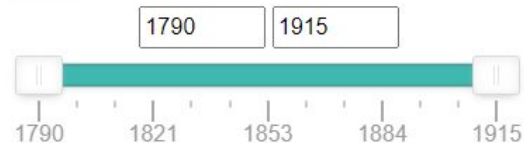
Below the text is a button labeled "Frise chronologique". The map background shows buildings and streets, with a search icon in the bottom left corner.

Ex : photo

Adresse

Ex : rivoli

Période



? Le filtre temporel permet de faire varier l'affichage des points préalablement chargés sur la carte sans lancer une nouvelle recherche. Données chargées pour la période **1790-1915**.

Localisation

Dessinez l'emprise de votre zone de recherche avec l'outil de dessin disponible sur la carte.

Lancer la recherche

Suivi individuel des commerces

ANNÉE – 1861

VUILLEMIN

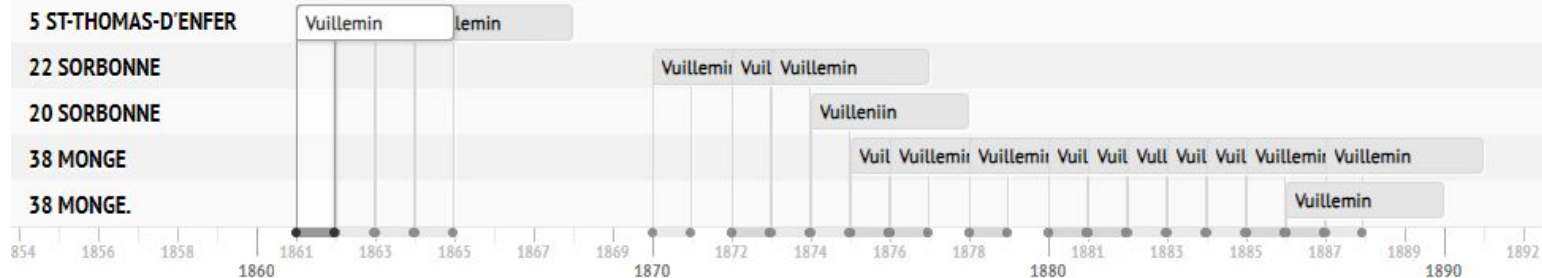
cartes géographiques

5 St-Thomas-d'Enfer

Source : DidotBottin_1861

Identifiant de l'entrée : 1d8338bc-8614-54cc-afd5-0d412ee29a35

Nombre de ressources liées : 27



Livrables disponibles



Publications *sur le thème “extraction depuis les annuaires”*

N. Abadie, E. Carlinet, J. Chazalon et B. Duménieu, “A Benchmark of Named Entity Recognition Approaches in Historical Documents Application to 19th Century French Directories”, in Proc. DAS 2022,

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

S. Tual, N. Abadie, E. Carlinet, J. Chazalon, et B. Duménieu, "A Benchmark of Nested NER Approaches in Historical Structured Documents", in Proc. ICDAR 2023, <https://github.com/soduco/paper-nestedner-icdar23-code/>

P. Bernet, J. Chazalon, E. Carlinet, A. Bourquelot et E. Puybareau, “Linear Object Detection in Document Images using Multiple Object Tracking”, in Proc. ICDAR 2023, <https://github.com/EPITAResearchLab/bernet.23.icdar>

S. Tual, N. Abadie, B. Duménieu, J. Chazalon et E. Carlinet. Création d’un graphe de connaissances géohistorique à partir d’annuaires du commerce parisien du 19ème siècle: application aux métiers de la photographie. IC 2023, 34èmes journées francophones d’Ingénierie des connaissances, Strasbourg, France, 3-7 July 2023. ⟨hal-04121643⟩.

https://github.com/soduco/ic_2023_photographes_parisiens

Données sur le thème “*extraction depuis les annuaires*”

A Dataset of French Trade Directories from the 19th Century (FTD)

DOI [10.5281/zenodo.6394464](https://doi.org/10.5281/zenodo.6394464)

A Dataset of French Trade Directories from the 19th Century for Nested NER task

DOI [10.5281/zenodo.7864174](https://doi.org/10.5281/zenodo.7864174)

Collection SoDUCo sur NAKALA : <https://nakala.fr/collection/10.34847/nkl.abe0gxah>

Annuaire historiques parisiens, 1798-1914. Extraction structurée et géolocalisée à l'adresse des listes nominatives par ordre alphabétique et par activité dans les volumes numérisés

<https://nakala.fr/10.34847/nkl.98eem49t>

Les datasets par professions incluant les données et les liens sont disponibles sur le SPARQL endpoint du projet : <https://dir.geohistoricaldata.org/>

Outils et modèles *sur le thème “extraction depuis les annuaires”*

Pré-traitements et segmentation du canevas : ouverture en cours, détection de segments disponible à <https://github.com/EPITAResearchLab/bernet.23.icdar>

OCR : benchmark DAS 2022 : <https://github.com/soduco/paper-ner-bench-das22>

NER : idem + modèles sur HuggingFace Hub, nouveaux prototypes en cours de finalisation

NER imbriqué (code, données et modèles) : <https://github.com/soduco/paper-nestedner-icdar23-code>

Interface Web de visualisation cartographique du graphe des photographes parisiens : https://soduco.geohistoricaldata.org/ic_2023_photographes_pariens/

Scripts d'extraction d'entrées pour une activité donnée, de liage, de structuration en RDF, de requêtes SPARQL types pour explorer les graphes géohistoriques et interface Web de visualisation cartographique capable d'intégrer plusieurs graphes:

https://soduco.geohistoricaldata.org/atelier_graphes_geohistoriques_annuaires/

Il reste quelques places pour l'atelier de demain après-midi !

Travaux en cours & futurs



Evaluation des étapes de la chaîne de traitement

Extraction d'entrées :

- échantillonnage → typologie des erreurs
- collection → détection de cas aberrants

Géocodage → détection de cas aberrants

Liage → évaluation de la cohérence intrinsèque, détection de liens erronés, etc.

Analyses en aval → détection de zones d'ombres

Objectif: **Aller vers une qualification plus fine des données au regard des problématiques géo-historiques visées.**

Amélioration de la chaîne d'extraction

- Modèles de langue hybrides texte / vision et analyse du texte en flux
- Nouveaux systèmes de transcription (cf présentation C. Kermorvant à suivre)

Liage en amélioration intensive

- Développement de chaînes de traitement, pour permettre l'extraction de données pour des activités spécifiques (semaine “graphe géo-historique” en oct. 23)
- Liage massif des données qui semble possible à court/moyen terme

Généralisation de l'approche

- **Ateliers : améliorer la réutilisabilité, la diffusion des données, outils et modèles, pistes de généralisation de la chaîne de traitements...**
- **Mezanno (plan quadriennal BnF) : outils pour constituer un corpus à partir de ressources IIF et l'annoter de façon semi-automatique, en toute autonomie**
- Plus de tests sur des corpus voisins : annuaires étrangers, de propriétaires, dictionnaires, débats parlementaires...

RÉSERVE

Distribution des outils

LIB



CLI

API REST

Batch

annotation assistée

I.2 - Des sources géographiques anciennes à l'analyse géohistorique

Reconnaissance

Classification

Géo-référencement

Structuration

Image → Géométries



Scan → Texte

Ravrio et comp., fabr. de bronzes et curiosités, r. Richelieu, 93; la fabrique rue Montmartre, 161.

Géométries → Annotations



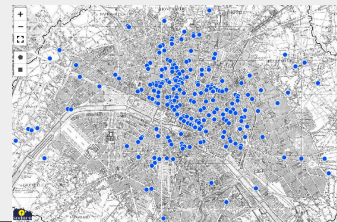
Texte → Entités nommées spatiales

Ravrio et comp. PESE , fabr. de bronzes et curiosités ACT
r. Richelieu LOC , 93 CARDINAL ; la fabrique FT
rue Montmartre LOC , 161 CARDINAL

Données géométriques → Données géographiques



Entités nommées spatiales → Données géographiques

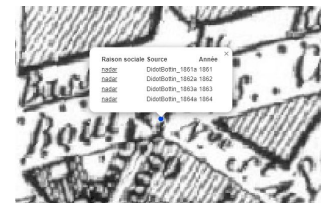


Données géographiques → Données géohistoriques

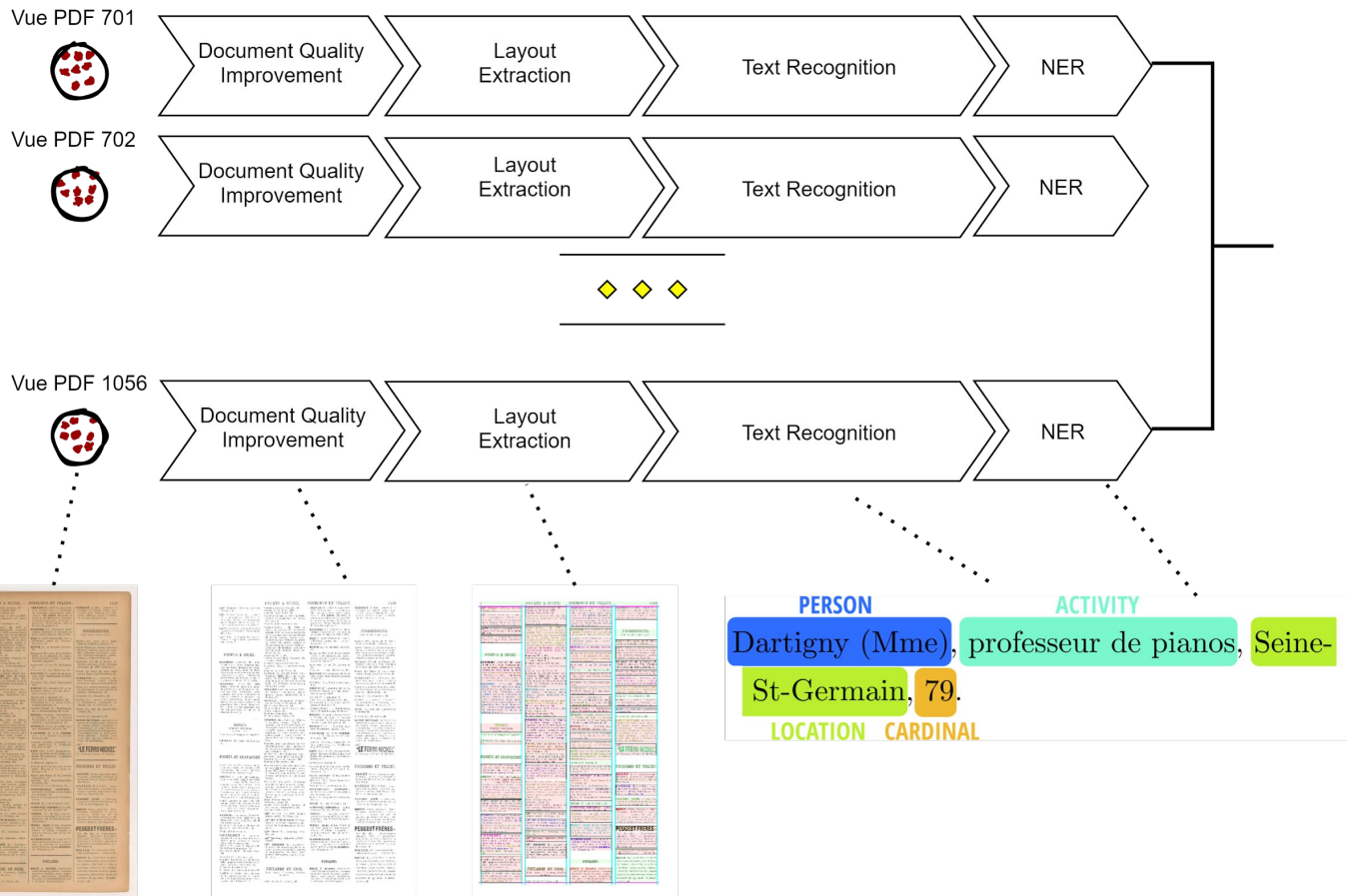


Aborder le liage ici ?

Données géographiques → Données géohistoriques



Chaîne de traitement automatique : de l'image aux infos sémantiques



Chaîne de traitement automatique : de l'image aux infos sémantiques

Vue PDF 701



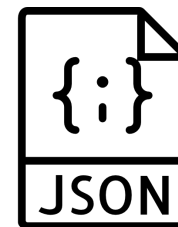
Vue PDF 702



Vue PDF 1056



un JSON par vue



V2 data (2022-06): ~10 M d'entrées



Chaîne de traitement automatique : ajout des infos spatiales

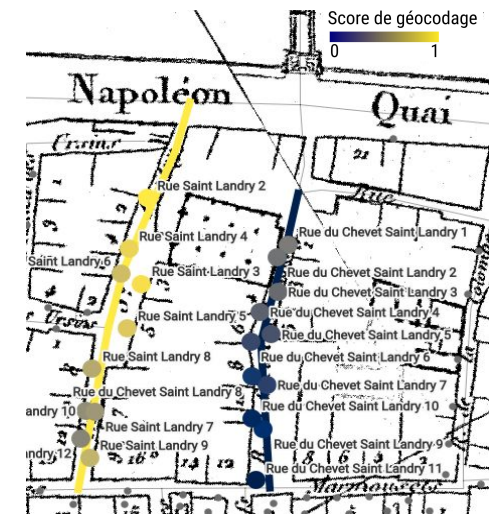
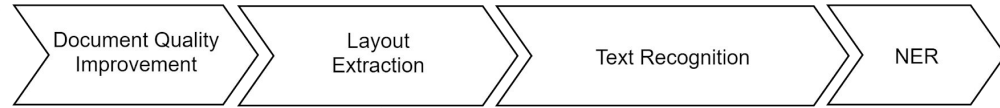
Vue PDF 701



Vue PDF 702



Vue PDF 1056



PERSON
Dartigny (Mme), professeur de pianos, Seine-St-Germain, 79.

ACTIVITY

LOCATION **CARDINAL**

Approche scientifique du projet SoDUCo

Analyse critique croisée de 2 collections

Atlas et annuaires

Exploitation des **redondances spatiales et temporelles**

Détection d'erreurs ou de changements, stabilisation de la reconnaissance

Extraction d'information **pragmatique**

Assister l'utilisateur expert, viser l'annotation collective dans une certaine mesure

Compromis entre **généricité** et **spécificité**

Viser une solution fonctionnelle sur notre corpus, avec des composants réutilisables

Free/open, accessible, indexé, reproductible

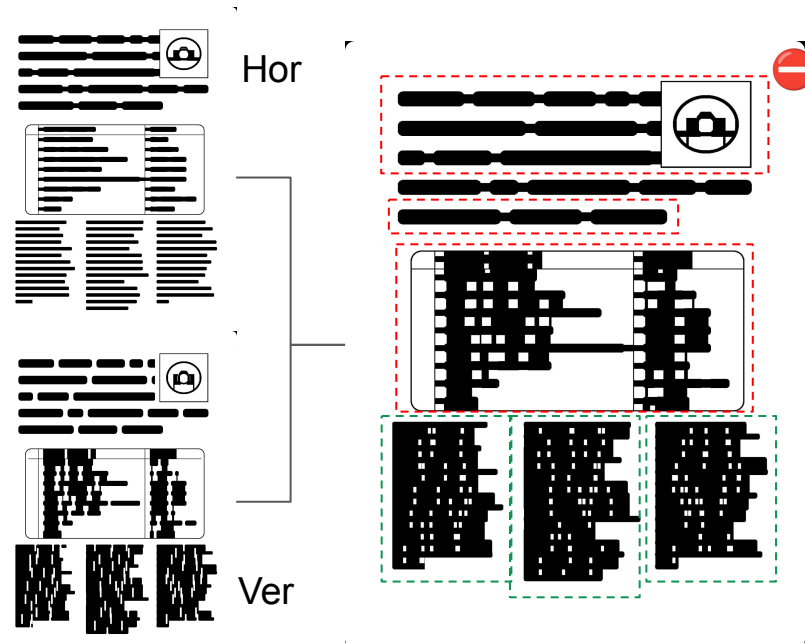
[Dépôts GitHub](#), <https://geohistoricaldata.org>, jeux de données, publications...

Segmentation du canvas (pour les blocs)

XY-Cut (Projection d'histogramme récursif)

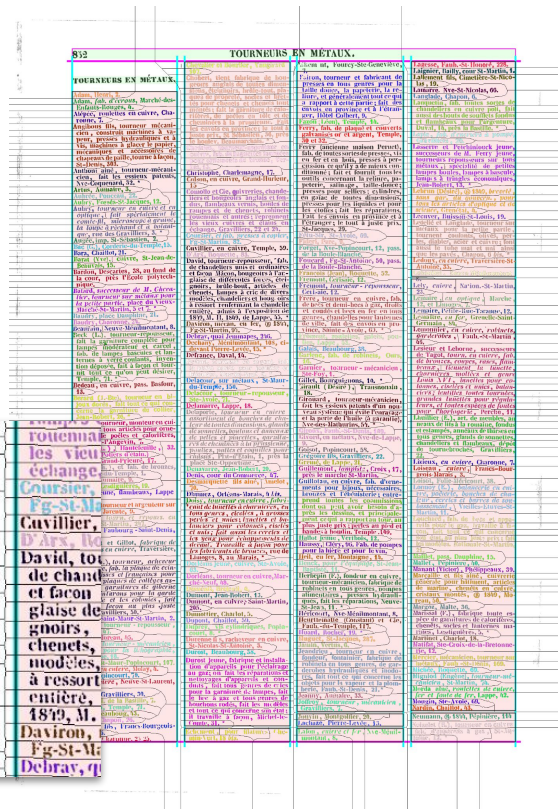
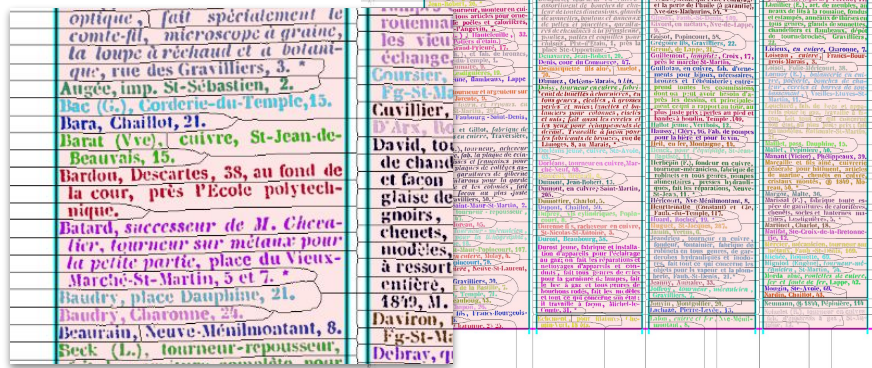
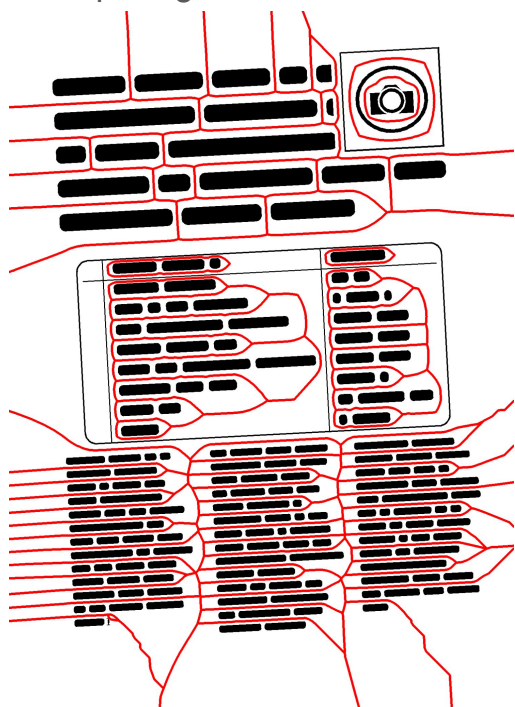


Smearing (ouverture/fermeture morphologique)



Segmentation du canvas (pour les lignes)

Ligne de partage des eaux



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.

Baridou, Descartes, 38.

Bariquand, 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 placage. Le bois on

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

province, 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, juseaux 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, 16.

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 eaux 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, 16.

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-

Phases d'ouverture des outils

Phase 0 : partage en **interne**, validation de l'utilité dans le cadre spécifique de SoDUCo, publication de prototypes isolés (recherche reproductible)

Phase 1 : tests avec d'**autres projets pilotes**, identification des éléments les plus utiles sans adaptation réelle (cas du projet AGODA au BnF Datalab)

Phase 2 : **adaptation à d'autres projets similaires** et choix d'une solution technique pour une évolution à long terme

Phase 3 : **ouverture complète** avec un soin particulier apporté à l'interopérabilité, la documentation, la maintenance, la sécurité, etc. en actant un positionnement et une démarche centrée sur quelques usages

Perspectives — Phase 2 d'ouverture en cours

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, intégration IIF...)
- **Partager** plus largement ces outils (*open source*, documentation, voire formation)

Liage en amélioration intensive

Développement de chaînes de traitement, pour permettre l'extraction de données pour des activités spécifiques (semaine “graphe géo-historique” en oct. 23)

Croisement massif des données qui semble possible à court/moyen terme

Master puis thèse de Solenn Tual

- Suivi des photographes
- Transposition au cadastre napoléonien

Amélioration des outils d'annotation et de visualisation

React App

Directory Viewer

Didot_1852a.pdf

700

EXPORT SAVE

Affichage

Entry

OCR NER

Liguez (P. PER), laminage de tous métaux pour bijouterie, orfèvres ACT, Chapon LOC, 18 CARDINAL ; fabrique de plaqué d'argent ACT, St-Martin LOC, 229 CARDINAL, anioien LOC 175 CARDINAL.

Comment

Tags

Checked

SAVE

LAMINEURS.

sale, Faub.-St-Martin, 69, No-
naindières, 12, de Cotte, 11, et
Penthière, 36. —

Hébert, en gros, Faub -St-Denis,
162.

Lefèvre fils aîné, Traversaire, 34

Piot et Lefèvre, laiterie en gros,
Amsterdam, 39. —

Poinsoi, nourrisseur-crémier, lait
d'ânes et de chèvre à domicile,
Chabrol, 82. —

Sarasin, Faub.-St-Martin, 270.

LAMINEURS.

Albaret aîné, fondeur et apprêteur
de métaux, fab. de maillechort,
prepara pour MM. les orfèvres,
lunetiers, couteliers, monteurs
de boîtes et garnisseurs, moulures
mat et vif, filets d'ébenistes,
entrepren la pièce de fonte, etc.,
place de l'Ancien-Marché-St-
Martin, 7. —

Bachollet et Cavillier, fab. de pla-
qué or et argent, laminages à fa-
çon de toute sorte de métaux,
St-Maur-Popincourt, 134. —

Cailar (J.-M.), fonderie, lami-
nage et tréfilerie de maillechort
pour orfèvrerie, chirurgie et
coutellerie, laminage à façon
de toutes sortes de métaux, exécuté

jaune et 119 rouge, Montgolfier
6, Marché-St-Martin. —

Liguez (P.), laminage de tous
métaux pour bijouterie, orfèvres,
Chapon, 18 ; fabrique de plaqué
d'argent, St-Martin, 229, anioien
175. —

Naudin (F.), rue Montmorency,
14.

Oeschger (L.), Mesdach et Cie,
fonderie, affûage de cuivre, de
zinc et de plomb, force hydrau-
lique de 80 chevaux, quatre lami-
noirs martinet, à Beaucroix-St-
Vaast, près Arras (Pas-de-Calais),
1849 ; maison à Paris, r. St-
Paul, 28. —

Pencé et Prévost, Bailly-St-Mar-
tin, 13.

Roussseau, ancienne maison Clic-
quot, 1849, outils de bijoutiers,
rouleaux, acier fondu, unis ou
gravés, filés, fil et plané en tous
genres, laminage de métaux, r.
Beaubourg, 50. —

Wilkens, gendre et successeur de
Soovens, apprêteur, découpeur
et estampeur pour MM. les bijou-
tiers, achète le vieux cuivre et
l'innaille, Graviillers, 24, et pass.
de Rome, 36.

LAMPES (fab. de). Voyez aussi
FERBLANTIERS-LAMPISTES.

soduco.geohistoricaldata.org

Favre_et_Duchesne_1798

Annotations

ANOTTAUX, rue Grenetat, n.° 46, —des Amis-de-la-Patr.

Basset, rue Bailleul, n.° 238, —des Gardes-Françaises.

Bazanne, rue Madeleine, n.° 1429, —du Roule. Berçot, rue Guénégaud

Bourgeois, rue Quincampoix, n.° 58, —des Lombards.

Boussode, rue Eloy, n.° 28, —de la Cité. Bouvrain, rue du Cimet-Je

Buffard, rue des Moineaux, n.° 413, —de la Buttes-des-Moul.

Collin, rue du Bouloy, n.° 21, —de la Halle-au-Bled. Ve Collin, rue de

Chatelain, rue Antoine, n.° 41, —des Droits-de-Homme.

Damour, rue du Rocher, n.° 518, —du Roule.

Delaporte, faub. Laurent, n.° 163, —du Nord. Double, rue de Norma

Ducrot, rue Germain, n.° 87, —du Muséum.

Duherche, rue du faub. Denis, n.° 30, —Poissonnière.

Duprez, rue Beauregard, n.° 224, —de Bonne-Nouvelle.

Durand, Carré-Martin, n.° 9, —des Graviillers. Durand, rue du Jour, n.

7 of 541 · Page 30

Page 27 Page 30 Page 31 Page 32 Page 34

AUBERGISTES.

Amis-de-la-Patr. —des Amis-de-la-Patr.

Basset, rue Bailleul, n.° 238, —des Gardes-Françaises.

Bourgeois, rue Madeleine, n.° 1429, —du Roule.

Boussode, rue Eloy, n.° 28, —de la Cité. Bouvrain, rue du Cimet-Je

Buffard, rue des Moineaux, n.° 413, —de la Buttes-des-Moul.

Collin, rue du Bouloy, n.° 21, —de la Halle-au-Bled. Ve Collin, rue de

Chatelain, rue Antoine, n.° 41, —des Droits-de-Homme.

Damour, rue du Rocher, n.° 518, —du Roule.

Delaporte, faub. Laurent, n.° 163, —du Nord. Double, rue de Norma

Ducrot, rue Germain, n.° 87, —du Muséum.

Duherche, rue du faub. Denis, n.° 30, —Poissonnière.

Duprez, rue Beauregard, n.° 224, —de Bonne-Nouvelle.

Durand, Carré-Martin, n.° 9, —des Graviillers. Durand, rue du Jour, n.

Assister le travail humain avec des outils automatiques

Explorer



Visualiser



Annoter



Analyser



API IIIF

E.g {BNF Gallica

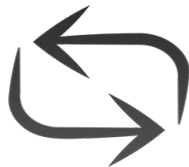
Export structuré

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)

Mezanno – plan quadriennal BnF

Outils de constitution de corpus, annotation, visualisation

Libres et ouverts, minimaux

Fonctionnalités visées : Autonomie des chercheurs SHS

- Pas d'installation requise, pas de serveur requis : une page à ouvrir
- Utilisation de ressources de calcul distantes, à la demande (OCR, NER...)
- Constituer son corpus
- Extraire des données de façon assistée
- Possibilités d'export

Démarrage courant 2024

Un géocodeur historique sensible aux temps valides

Adresse
référence spatiale indirecte
r. St. Landry Cité, 2

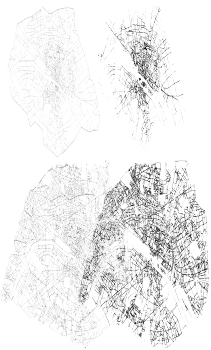
Année cible
1838




Géocodeur



Quadruplets
rues & adresses



-  **Analyse syntaxique (parsing)**

< rue >	< quartier >	< numéro >
r. St. Landry	Cité,	2
-  **Recherche documentaire (collect, filter & score)**
-  **Tri temporel (reranking)**

