

Compter les mots pour remonter le temps : Gallicagram et Gallicagrapher, deux outils d'exploration des archives numérisées de la BnF

Les archives numérisées de la Bibliothèque nationale de France, Gallica, constituent un trésor de données ouvertes. Les auteurs ont développé deux logiciels de lexicométrie, qui mesurent et permettent de visualiser l'évolution de l'usage des mots au cours du temps, et facilitent l'accès au contexte des occurrences.

BENJAMIN AZOULAY

Administrateur-élève des affaires maritimes, ministère de la Mer et École normale supérieure Paris-Saclay

BENOÎT DE COURSON

Doctorant au Max Planck Institute de Freiburg (Allemagne)

WILL GLEASON

Développeur indépendant

La Bibliothèque nationale de France (BnF) dispose d'un des plus riches fonds imprimés au monde. Elle est en effet dépositaire du dépôt légal : chaque livre ou numéro de presse publié depuis 1537 doit y être déposé pour archivage et désormais pour consultation par les chercheurs.

Depuis 1997, la BnF entreprend une numérisation massive de ses fonds et les verse en libre accès sur la plateforme Gallica. Cette « bibliothèque virtuelle de l'honnête homme¹ » est aussitôt devenue un outil de travail incontournable pour les chercheurs en humanités.

De par son volume (400 000 livres et 4,5 millions de numéros de presse ocrésisés² en français), Gallica se prête à merveille à la lexicométrie, c'est-à-dire au traitement quantitatif des textes. L'accès transparent aux données est ici un atout précieux. Développés à cette fin, les sites Gallicagram et Gallicagrapher visent à offrir aux chercheurs en humanités une interface Web pour exploiter et visualiser ces données.

Les outils lexicométriques existants n'appliquent guère les principes de la « science ouverte ». L'outil de référence en la matière, Google Books Ngram Viewer³, est peu utilisé par les chercheurs, embarrassés par la constitution opaque de son corpus et déçus de ne pouvoir accéder au contexte des occurrences⁴. Plus souple et transparent, Frantext⁵ permet des relevés syntaxiques dans près de 6 000 œuvres (soit 266 millions de mots), mais son volume est insuffisant pour des analyses quantitatives diachroniques⁶, et Frantext (par ailleurs payant) a récemment cessé d'assortir les données récoltées de graphiques.

Gallicagram et Gallicagrapher : deux logiciels pour démocratiser le traitement automatique de Gallica

Gallicagram⁷ est un logiciel permettant de visualiser l'évolution de l'usage des mots au cours du temps, en fouillant, parmi d'autres corpus, la presse et les livres numérisés de Gallica. La croissance des données interrogées à partir de la Révolution rend le corpus de presse (certainement le plus intéressant pour l'historien) particulièrement fiable entre 1789 et 1950.

Disposant d'un corpus ouvert, l'historien est à même de savoir dans quoi il cherche, et d'accéder au contexte des occurrences. Pour ce faire, Gallicagrapher⁸ exploite les API⁹ de Gallica et présente le contexte immédiat de chaque occurrence, directement dans le logiciel, sur le modèle de Frantext. Cela facilite l'analyse des résultats correspondant aux courbes affichées, ce qui permet, par exemple, de lever les ambiguïtés sur les homonymes et les erreurs de reconnaissance optique des caractères.

Les deux logiciels cherchent à appliquer les principes de la science ouverte et collaborative. Ils sont

1. <https://gallica.bnf.fr/edit/und/a-propos>

2. Le terme « ocrésisation » dérive de l'abréviation OCR : *Optical Character Recognition*, c'est-à-dire en français : « Reconnaissance optique des caractères » (ROC, peu utilisé). Techniquement, il s'agit du traitement d'une image (le texte est scanné, comme par une photocopieuse) sur laquelle on fait intervenir un logiciel de reconnaissance de caractères : le logiciel déchiffre les formes et les traduit en lettres.

3. <https://books.google.com/ngrams>

4. François Héran, « Les mots de la démographie des origines à nos jours : une exploration numérique », *Population*, vol. 70, 2015, p. 525-566.

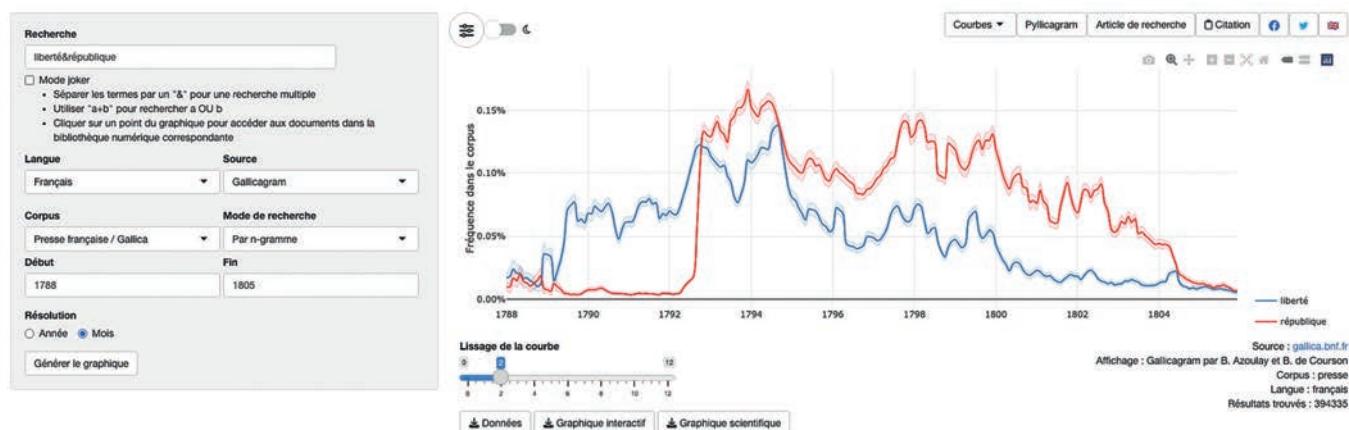
5. <https://www.frantext.fr>

6. Nous développons ce point dans l'article de Benoît de Courson, Benjamin Azoulay, Clara de Courson, Laurent Vanni et Étienne Brunet, « Gallicagram : les archives de presse sous les rotatives de la statistique textuelle », *Corpus*, n° 24, 2023. <https://doi.org/10.4000/corpus.7944>

7. <https://shiny.ens-paris-saclay.fr/app/gallicagram>

8. <https://www.gallicagrapher.com/>

9. Une API (*Application Programming Interface* ou « interface de programmation d'application ») est une interface logicielle qui permet de « connecter » un logiciel ou un service à un autre logiciel ou service afin d'échanger des données et des fonctionnalités.



Interface de Gallicagram : <https://shiny.ens-paris-saclay.fr/app/gallicagram>

Page 1 de 645 > >>

Document	Gallica	Contexte gauche	Pivot	Contexte droit
> (5) <i>Le Mémorial des Pyrénées : politk</i> 1888-08-07	Image	INFORMATIONS UES	GRÈVES	A PARIS Les terrassiers La réunion quotidienne des ouvriers grévistes terrassier
> (7) <i>Le Petit républicain : journal quoti</i> 1888-08-03	Image	donner leur adhésion à la	grève	des terrassiers. les cochers de fiacre Après deux discours des citoyens Winant e
> (4) <i>Les Chantiers de l'Exposition univ</i> 1888-08-01	Image	La	Grève	et l'Exposition de 1889. — Les Chantiers de Paris. — Les Travaux de Paris. — Une
> (12) <i>Le Cri du peuple : journal politiq.</i> 1888-08-01	Image	LA	GRÈVE	DES MINEURS (De notre correspondant) Saint-Etienne, 30
> (1) <i>La Lanterne de Boquillon / par A. I</i> 1888-08-19	Image	mirent en	grève	et les garçons de café «c extras », battus par Lozé II à Aboukir, se jurèrent de
> (7) <i>Le Clairon du Lot : journal monarc</i> 1888-08-07	Image	que les progrès faits par la	grève	dans d'autres corporations ne pouvaient qu'encourager les terrassiers à persévér
> (8) <i>Le Cri du peuple : journal politiq.</i> 1888-08-21	Image	La	grève	doit se transformer, comme s'est transformée l'industrie et la guerre elle-même
> (9) <i>Le Cri du peuple : journal politiq.</i> 1888-08-02	Image		GRÈVE	ET SYNDICAT Huit mille ouvriers se sont mis en
> (10) <i>Le Radical algérien : paraît tous l</i> 1888-08-05	Image	E» c est comme ce a partout 1 V La	grève	des mineurs. — On écrit du Saint Etienne
> (7) <i>Le Cri du peuple : journal politiq.</i> 1888-08-03	Image	UNE NOUVELLE	GRÈVE	Les verriers à vitres de Saint-Etienne Le verre à vitre. — L'usine Velin. — Un r
> (8) <i>La Presse</i> 1888-08-02	Image	Une	grève	se résume toujours pour les travailleurs par ces trois phases: Plus de salaires,
> (9) <i>L'intransigeant</i> 1888-08-06	Image	Le jour où toutes les corporations du bâtiment se mettraient simultanément ^n	grève	, on verrait ce que pèse le capital devant le travail, et il faudrait bien que l
> (8) <i>Le Radical</i> 1888-08-09	Image	LE DEVOIR DES PATRONS La	grève	des terrassiers semble entrer enfin dans la voie de l'arrangement
> (10) <i>Le Cri du peuple : journal politiq.</i> 1888-08-11	Image	La	grève	des ouvriers verriers. — Im- portante réunion ouvrière. — Pas de désordres
> (7) <i>Le Cri du peuple : journal politiq.</i> 1888-08-14	Image	Les terrassiers en	grève	ne font que réclamer de leurs exploiters les salaires et les conditions que le

*Open Source*¹⁰, ce qui permet le réemploi du code dans le cadre de projets tiers. Gallicagram est aussi *Open Data* : la base de données constituée par le décompte des milliards de mots des corpus numérisés de Gallica (presse et livres français) est accessible par API¹¹. Les deux logiciels collaborent puissamment à travers leurs API respectives : Gallicagram fournit ses données au graphique affiché par Gallicagraphe, qui lui renvoie le contexte des occurrences (figure ci-dessus).

Ceux-ci étant destinés à une population de chercheurs inégalement à l'aise avec l'informatique, l'ergonomie est un enjeu central. Gallicagraphe ne fait au fond qu'enrober les API de Gallica pour présenter de façon plus intuitive et plus maniable des données que le chercheur aurait pu trouver manuellement.

Mais il y a fort à parier que sans ce tour de force ergonomique, il n'en aurait pas eu le courage. Notons aussi que le développement de telles applications Web interactives a été facilité par les avancées récentes des langages de programmation (respectivement Shiny et React pour les deux logiciels).

Des obstacles persistants à l'ouverture des données

Ces logiciels se sont développés en tirant profit des données ouvertes de Gallica – mais aussi, avouons-le, en contournant les barrières à l'entrée de sources qui le sont moins.

Nous nous sommes surtout heurtés à Retronews, un service développé par BnF-Partenariats pour

Fonction « contexte » de Gallicagraphe : <https://www.gallicagraphe.com/>
Gallicagram utilise l'API de Gallicagraphe pour proposer une présentation analogue sous ses graphiques.

10. https://github.com/regicid/docker_gallicagram, <https://github.com/gleasonw/gallica-grapher>

11. <https://github.com/regicid/pylicigram>