

Emmanuelle Bermès 18/10/2017

LE PROJET DSR-CORPUS

- Un projet inscrit au plan quadriennal de la recherche de la BnF 2016-2019
- Objectifs:
 - préfigurer un service de fourniture de corpus numériques à destination de la recherche
 - fournir à des chercheurs des données et des outils pour les analyser, dans le respect du droit d'auteur et de la vie privée.
- 3 années d'expérimentation (archives web, numérisation, métadonnées) + 1 année de bilan
- → Fiche dans la base ANIR : http://c.bnf.fr/fom

QUELS CORPUS?



Gallica + Gallica intra muros 4.5 millions de documents numériques http://gallica.bnf.fr



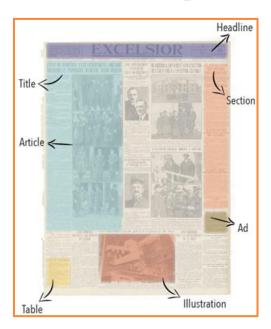
Archives de l'Internet 793 To

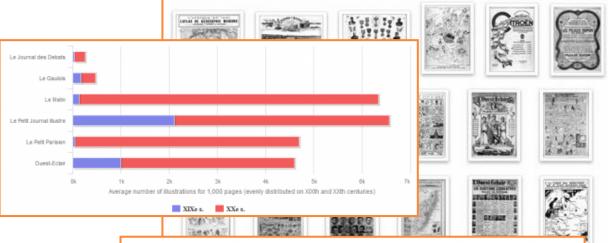
Métadonnées + de 20 millions de notices http://data.bnf.fr

AVANT CORPUS... (LE CORPUS COMME SOURCE)

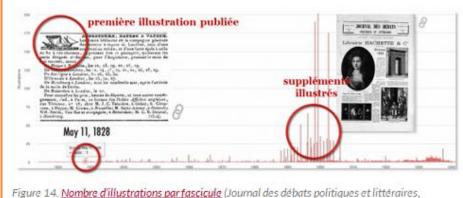
1814-1944)

The Europeana Newspapers project



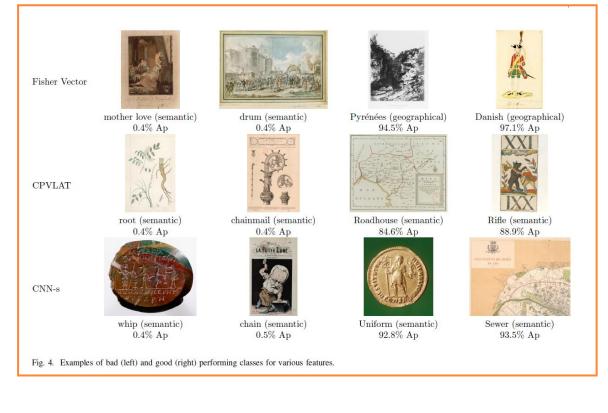


JP Moreux, "Approches innovantes pour la presse ancienne numérisée : fouille et visualisation de données" in *Carnet de recherche à la bibliothèque nationale de France*, 3 décembre 2016 https://bnf.hypotheses.org/208



AVANT CORPUS... (LE CORPUS COMME BAC À SABLE)

• ETIS : le projet ASAP

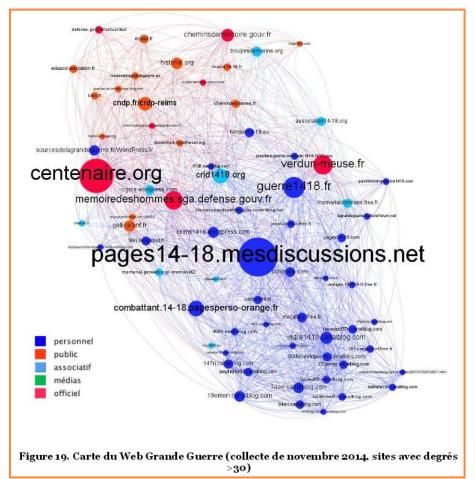


David Picard, Philippe-Henri Gosselin, Marie-Claude Gaspard. "Challenges in Content-Based Image Indexing of Cultural Heritage Collections." *IEEE Signal Processing Magazine*, Institute of Electrical and Electronics Engineers, 2015, 32 (4), pp.95 – 102

Corpus from http://images.bnf.fr

AVANT CORPUS... (LE CORPUS COMME INTERFACE)

• Project « le devenir en ligne du patrimoine numérisé : l'exemple de la Grande Guerre » - Labex « les passés dans le présent »



Valérie Baudouin, Zeynnep Pelhivan : Cartographie de la Grande Guerre sur le Web : Rapport final de la phase 2 du projet "Le devenir en ligne du patrimoine numérisé : l'exemple de la Grande Guerre"

https://hal.archivesouvertes.fr/hal-01425600

BILAN ANNÉE 1

Partenariat avec l'équipe Web90 (CNRS/ISCC)



cherche avan	cée (2)				
Mot(s):	mintel			Rechercher un ou plusieurs mots ou une "expression exacte"	
Proximité :	Mort.	Not2	\$6.	Rechercher des mots qui apparaissent à une distance de N	
	Au moirs un de ces deux champs doit être renseigné pour lancer une recherche			11.08.000000000000000000000000000000000	
Exclure le(s) mot(s) : (facultatif)				Exclure les documents où ce(s) mot(s) apparaissent	
res (facultatifs)					
Date :				Limiter la recherche à une date (aaaa, mm/aaaa ou jimm/aaaa)	
Ou période :	de 1997	a 1009		Limiter la-recherche à une période (assa, mm/assa ou jimm/assa)	
Hôte, nom de domaine, extension :	h			Limiter la recherche à un hôte (gallica bril ft), un nom de domaine (bril ft) ou une extension (ft)	
Format :				Préciser le format de fichier des documents (témi, text, pdf, image, audio, video, word, excel, powespoint, other)	

d'archives de l'internet

Les archives de l'internet comme sources : méthodes et représentations

ALABY AND RECREDIES SUBLES COLLECTIONS

IL ÉTAIT UNE FOIS DANS LE
WEB. 20 ANS D'ARCHIVES DE
L'INTERNET EN FRANCE

Valérie BEAUDOUIN, en
dans le présent), Sophie
MUSIANI, enseignante-ch
chercheuse (ISCC / Webs



fétent cette année leurs 20 ans, et la oil sur le dépôt légal du veb ses 10 ans. Le colloque « Il était une fois dans le veb, 20 ans d'archives de l'internet en Fance » organis bejar la Bibliothèque nationale de l'audiovisuel et l'institut national de l'audiovisuel avec le concours de l'Engliège du protet ANR Web20; se tiendra la 23 an vermbre 2016. Il retracerar l'histoire de la préservation de ce patrimoine simulier. Valérie BEAUDOUIN, enseignante-chercheuse (Télécom-ParisTech / Labex Les passés dans le présent), Sophie GEBEIL, enseignante-chercheuse (Unix. Aix-Marseille), Francesca MUSIANI, enseignante-chercheuse (ISCC / Web90), Valérie SCHAFER, enseignante-chercheuse (ISCC / Web90), Marie-Luce VIAUD, cheffe de projet recherche et développement (Ina), Dana DIMINESCU, enseignante-chercheuse (Télécom-ParisTech)



Durée : 59 min

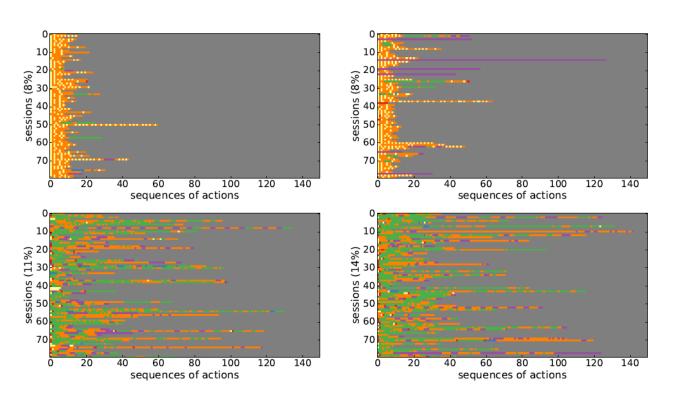


Video recordings: http://c.bnf.fr/fse
Wohlog: http://wohcornors.hymothe

Weblog: http://webcorpora.hypotheses.org/

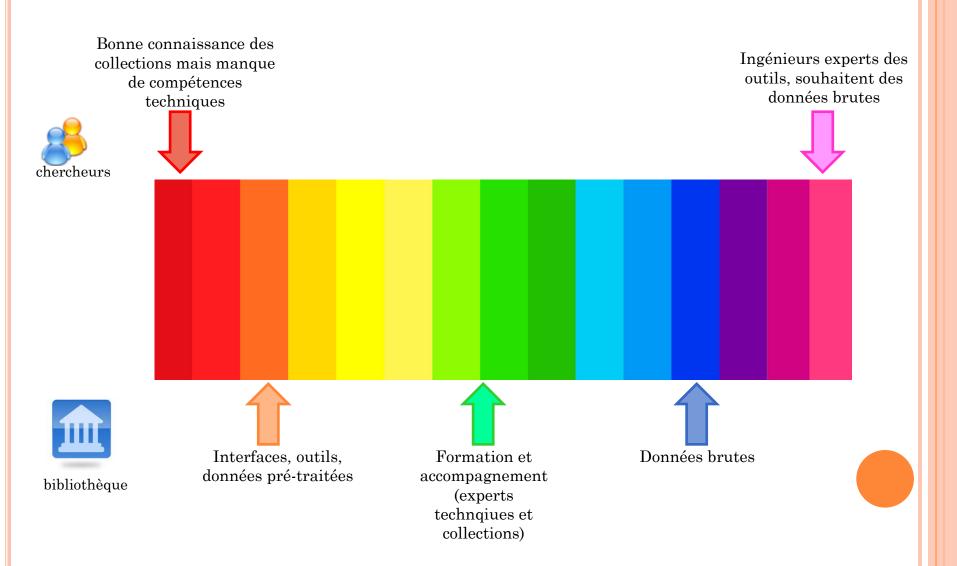
Bilan année 1 – Projets connexes

o Bibli-Lab : analyse des logs de Gallica



Etude réalisée par Adrien Nouvellet, Télécom Paristech

Une grande diversité de situations...



CONSTRUIRE L'AVENIR?

Accès distant

Infrastructure informatique

Espace physique



Aspects juridiques Offrir aux chercheurs, dans les emprises de la Bibliothèque, des outils de fouille et d'exploration de textes et de données sur des corpus numériques de la BnF

 Proposer des environnements scientifiques et techniques (plate-forme sécurisée, logiciels, assistance d'experts...) pour explorer, dans le respect des dispositions réglementaires, les corpus numériques de la BnF

> CONTRAT D'OBJECTIFS ET DE PERFORMANCE 2017-2021



Aspects organisationnels

Compétences et savoir-faire Ressources humaines

Vers un laboratoire d'étude et d'analyse de corpus numériques à la BnF



 $\underline{emmanuelle.bermes@bnf.fr}$



@figoblog