#### L'archivage du web français à l'ère de l'IA

Séminaire du weblab du 27 novembre 2025



Géraldine Camile

(BnF DataLab, Département de la découverte des collections et de l'accompagnement à la recherche)

# Plan

- Panorama des collections des archives du web de la BnF
- Modalités de constitution des collections des archives du web à la BnF
- Les outils de consultation et de recherche de la BnF
- Identifier les services offerts par la BnF
- Impact de l'IA sur les archives du web, de la sélection à la recherche

Panorama des collection à la BnF



# Fear of a digital dark age La peur d'un web amnésique

Brewster Kahle lance la fondation Internet Archive en 1996

Les12 membres fondateurs d'IIPC signe l'acte de naissance du consortium en 2003 à la BnF.

https://web.archive.org/ https://netpreserve.org/

Des archives du web qui remontent à 1996, avec rachat de collections à Internet Archive et expérimentations de collecte

# La mission de dépôt légal de l'internet

La loi DADVSI (2006) et son décret d'application (2011), intégrés au Code du patrimoine, instituent le « dépôt légal du web ».

Dans la continuité des missions des institutions patrimoniales : « Rassembler, conserver le patrimoine culturel français » et y donner accès.

Deux établissements dépositaires : INA (sites de radio et télévision), BnF (le reste du web français).

Les contenus restent soumis au droit d'auteur et les données collectées sont des données patrimoniales.

Des collections qui illustrent ce qu'est le web, dans toute sa diversité de formes et de contenus, à une époque donnée





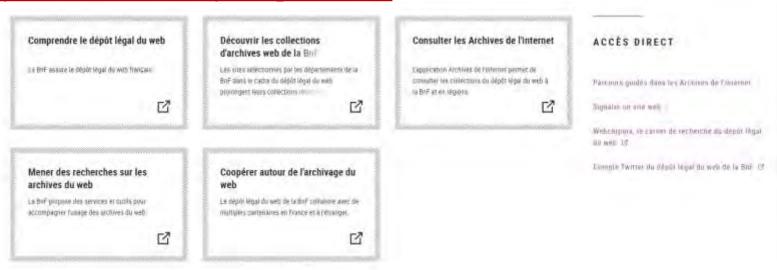




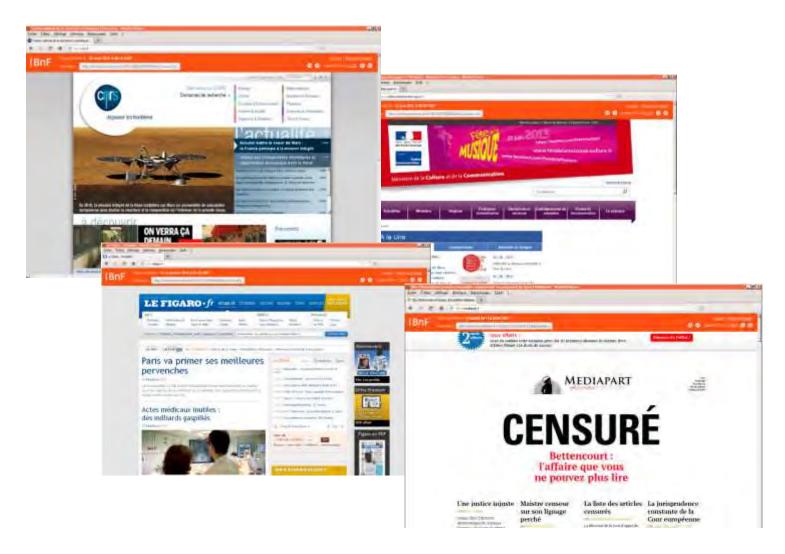




#### https://www.bnf.fr/fr/depot-legal-du-web



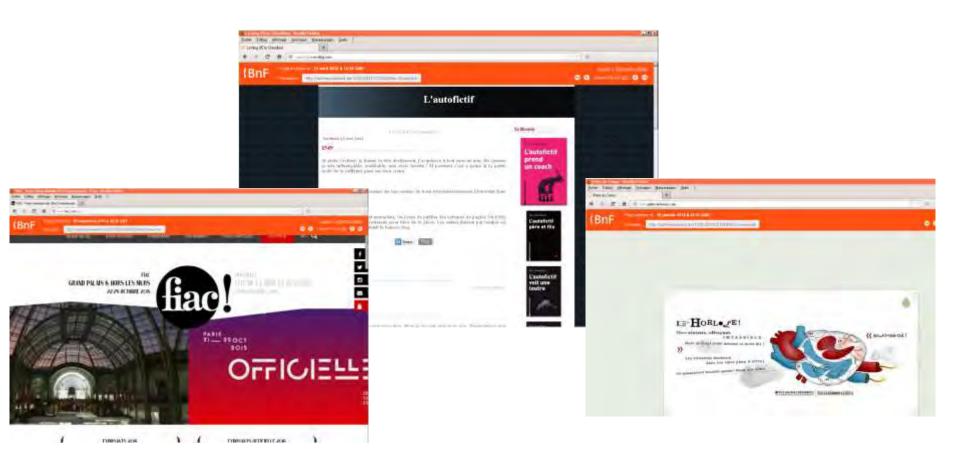
Des sites officiels, de référence, de médias



Des journaux en PDF, des livres numériques en EPUB



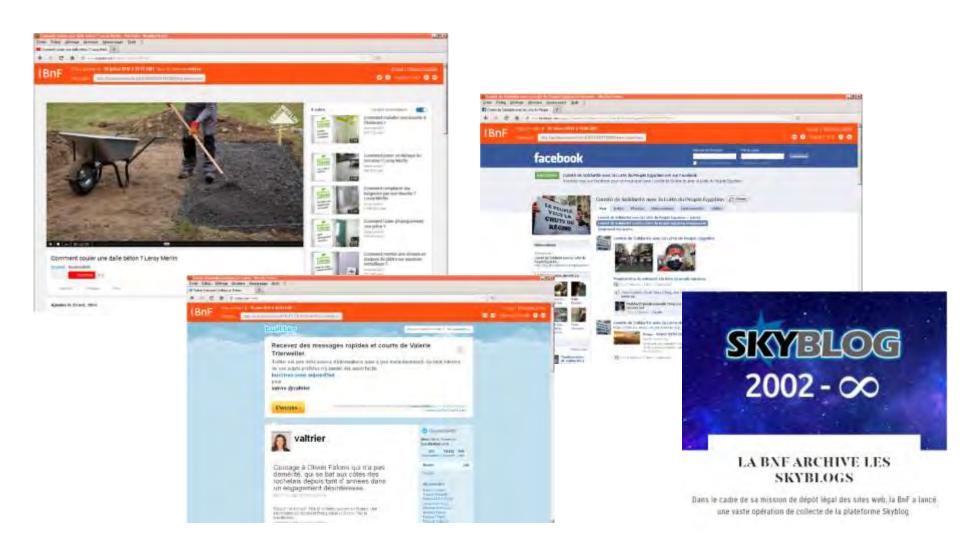
### Des sites sur l'art, la littérature en ligne



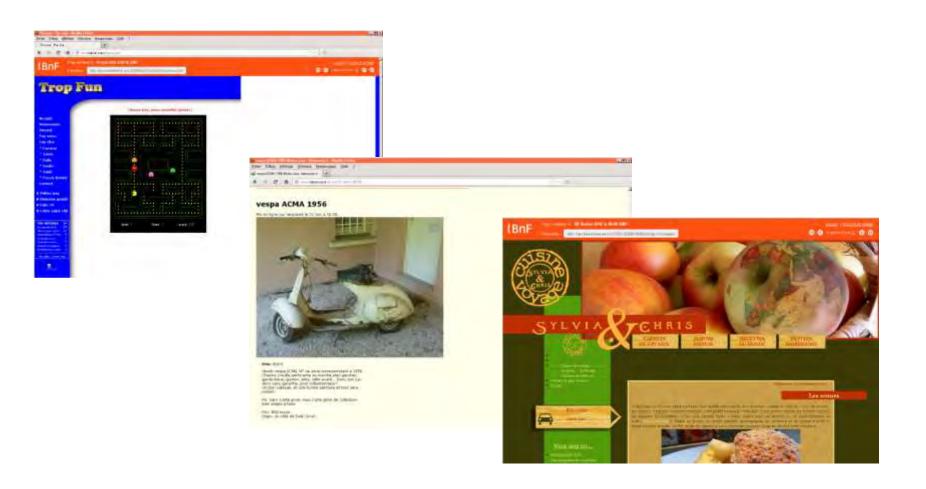
#### Des sites liés à un événement



#### Réseaux sociaux, vidéo et audio



#### ... témoins de la diversité du web



# des contenus qui évoluent dans le temps

http://solidarites-sante.gouv.fr/





March | Mileson product

Février 2020 Mai 2020

#### Qu'est-ce qu'une archive web?



Niels Brügger

- La notion de *reborn-digital material* a été developpée par Niels Brügger pour qualifier l'archive web.
- « Le web archivé est fondamentalement différent de tous les autres types de médias existants, y compris le papier, la radio, la télévision et même le web vivant ».
- Il n'y a pas de copie de la source comme dans un acte de photographie, mais une **transformation** en archive.
- L'archive est le résultat des **modalités techniques** de constitution des collections et de la configuration des outils d'accès.

=> Comprendre l'archive, c'est comprendre les collectes.

Brügger, N. (2016). Digital Humanities in the 21st Century: Digital Material as a Driving Force. *Digital Humanities Quarterly*, 10(3). http://www.digitalhumanities.org/dhg/vol/10/3/000256/000256.html

Modalités de constitution des collection à la BnF



#### Les collections du Dépôt légal de l'internet

- Des captures de sites web, **restitués** à partir des fichiers collectés tels qu'ils étaient au moment des **crawls**
- Dans la continuité des autres collections de la Bibliothèque
- Dépôt légal : tous les types de site et toutes les formes de publication
- On retrouve des équivalents des publications sur support, et aussi de nouvelles formes spécifiques au web

Deux modalités principales et complémentaires de constitution des collections

#### Collecte large

• Des échantillons de tous les sites du domaine français, une fois par an (listes fournies par l'AFNIC, OVH, etc.)

#### Collectes ciblées

- Dans la continuité des collections de la BnF
- Importance croissante des collectes projet : élections, COVID-19, actualité, enjeux environnementaux, intelligence artificielle, etc.
- Sélection qui obéit aux principes fondamentaux du dépôt legal : pas de jugement scientifique, moral, ni esthétique

# (BnF BnF Collecte du web

#### Voir les collectes



#### FICHE DE SITE 97 357

Collecte: Mouvements sociaux (MOUV)

Responsable : Sophie GEBEIL

Création : Sophie GEBEIL, le 29/05/2015 Mise à jour : Sarah TOURNERIE, le 24/11/2023

#### Paramètres de collecte

Etat : Inoctif

URL de départ : http://leschibanis.free.fr

URL supplémentaires :.

Type de collecte : cibiée

Budget : petit

Fréquence : 1 fois par an

Profondeur : hôte

#### Description du site

Thème: memoires et patrimoines de l'immigration

Mots clés : chibanis / immigrés âgés / travailleurs étrangers maghrébins / photographie

Notes de contenu : Ph Flash (SG)

Site constitué de photographies sans commentaires.

Notes techniques: DLN - 20170628 - mise en collecte d'urgence pour parcours guidé (transfert MOUV -> DDL); rajout URL supplémentaire (U profondeur domaine.

 - après collecte remettre en MOUV, responsable Sophie Gebeil, petit - 1 fois par an - hôte ; thème mémoires maghrébines DLN - 20170711 - Transfert de DDL (demandes externes) à MOUV

DLN - 20190529 - suppression de l'URL http://www.chibanis.com (remplacée par l'aricienne URL supplémentaire)

PHS - 20221006 - 35 photos actuellement. Voir si le site évolue, et si des photos sont ajoutées d'ici un an, sinon l'inactiver.

PHS - 20232411 - sité n'existe plus Désactivé le 24/11/2023 ST

#### Outil de curation

- Vue documentaire par collecte
- Saisie des paramètres de collecte
- Thème, mot clé, description
- Exemple:

leschibanis.free.fr fiche créée par Sophie Gebeil, en 2015, inactif depuis 2023

# {BnF

A défaut d'exhaustivité, un objectif de représentativité et d'échantillonnage

Les collectes ciblées :

quels principes de sélection ? Suspension du jugement esthétique, moral, scientifique : culture populaire comme culture savante, *mainstream* ou longue traîne, tout est collectable

Capturer les actions, les savoirs, mais aussi les idées, les représentations qui circulent sur le web sur un sujet donné : rendre compte de la diversité du web, media et espace d'échange



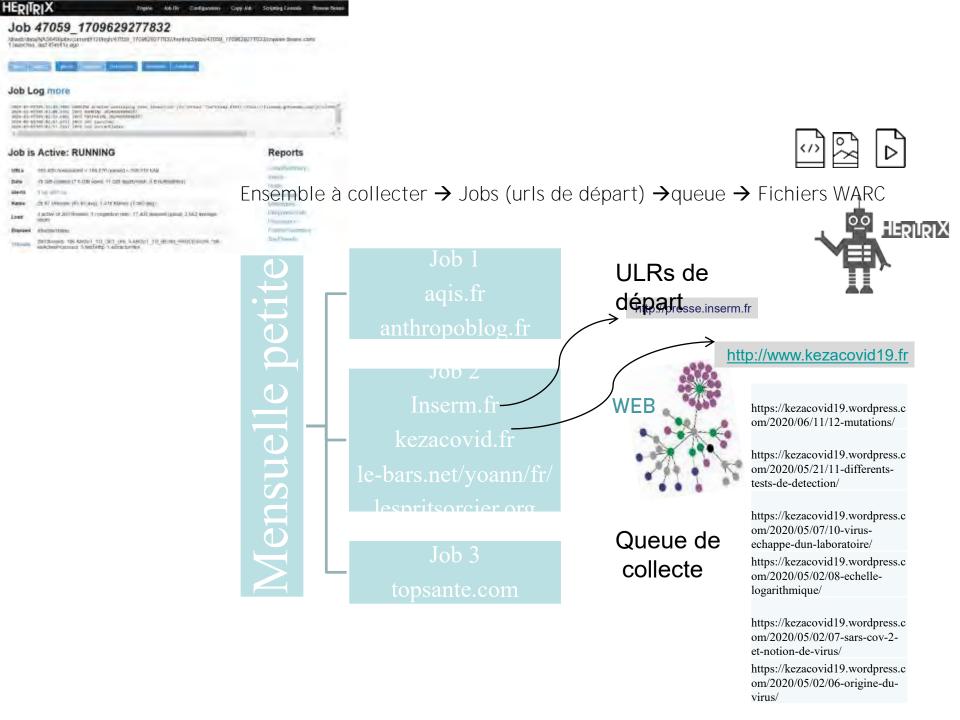
#### Le robot de collecte Heritrix

```
<html lang="fr" dir="ltr" prefix="content: http://purl.org/rss/1.0/modules/content/ dc: http://purl.org/</pre>
    <meta charset="utf-8" />
k rel="canonical" href="https://www.bnf.fr/fr" />
k rel="shortlink" href="https://www.bnf.fr/fr" />
<meta name="description" content="." />
<meta property="og:site name" content="BnF - Site institutionnel" />
<mata property="og:type" content="article" />
<meta property="og:url" content="https://www.bnf.fr/fr/node" />
<meta property="og:title" content="En ce moment | BnF - Site institutionnel" />
(meta name="Generator" content="Drupal 8 (https://www.drupal.org)" />
<meta mame="MobileOptimized" content="width" />
<meta name="HandheldFriendly" content="true" />
<meta name="viewport" content="width=device-width, initial-scale=1.0, shrink-to-fit=no" />
<meta http-equiv="x-ua-compatible" content="ie=edge" />
k rel="shortcut icon" href="/themes/custom/bnfsi/favicon.ico" type="image/vnd.microsoft.icon" />
<script>window.a2a config=window.a2a config||{};a2a config.callbacks=[];a2a config.overlays=[];a2a config.
a2a config.num services = 6;
a2a config.prioritize = ["facebook", "twitter", "email", "whatsapp", "linkedin", "pinterest"];
var a2a config = a2a config | { };
a2a config.onclick = 1;</script>
    <title>Accueil | BnF - Site institutionnel</title>
    «link rel="stylesheet" media="all" href="/sites/default/files/css/css WKsTfRhj5dffmDmprU3p5r7TI-KRK13U
style>.dropdown-language-item ,language-link.inactive-language {color:#c8cccf} /*Différencier le faux
.video is .vis-tech {max-height: 85vh} /*Limiter la hauteur des vidéos verticales*/
.glossaire h2 {color: white; background: #333; text-align: center}
dd {margin:0 0 1em 2em}
dt {color:black}</style>
<!--[if Ite IE 8]>
<script src="/sites/default/files/is/is VtafiXmRvoUgAzqzYTA3Wrjkx9wcWhipQ64ZnnqRamA.js"></script>
<script src="//tag.aticdn.net/18798/smarttag.js"></script>
<script src="https://tarteaucitron.io/load.is?domain=bnf.fr/fr&amp;uuid=91ba7c974c752a888f0a2765181efcbea3</p>
  «body class="view-page view--frontpage view--frontpage--page-1 path-frontpage has-glyphicons">
    (a href="#main-content" class="visually-hidden focusable skip-link">
      Aller au contenu principal
      (div class="dialog-off-canvas-main-canvas" data-off-canvas-main-canvas>
    <div class="left bar fixed">
       <div class="region region-left-bar-fixed">
    <div id="block-menuvertical" class="block block-block-content block-block-contentd7b2e8a9-7ccd-4c32-b6</p>
```

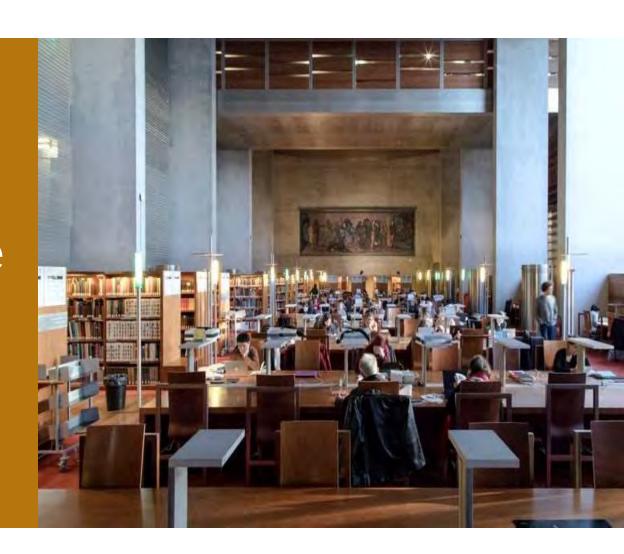
Le robot part d'une liste d'URL de départ et exploite les liens dans les codes sources des sites.

Des paramètres de collecte sont ajoutés. Certains paramètres peuvent être modifiés en cours de collecte.

Le robot copie les fichiers dans des fichiers de sauvegarde WARC, format standard pour la conservation.



Les outils de consultation





# Accès à la BnF et dans 22 bibliothèques partenaires

réseau des bibliothèques de dépôt légal imprimeur (BDLI)



# (BnF Les Archives de l'internet, une wayback machine





#### Rien n'indique:

- les origines techniques ou documentaires de la collecte
- ni la qualité de l'archive

Par exemple : Les 14 résultats présents ne sont pas pertinents car la vraie page est : http://leschibanis.free.fr/Menu1.html

#### Petite philologie du web

- Des sauts temporels lors de la navigation
- Une « page web » recomposée à partir d'éléments crawlés à différents moments
- Des pages inventées par l'extracteur Javascript du robot
- Des sites collectés à des fréquences différentes
- Des sites collectés profondément et d'autres superficiellement
- Des blocages du robot, des codes réponses http
- Des contenus « embed » qui laissent des traces fantômes
- L'absence de statistiques d'usage...

Les Archives de l'internet restituent ce que le robot a collecté....

#### Citer une « Archive de l'internet »

- Permalien
  - Même construction pour Internet Archives que pour la Wayback Machine d'Internet Archives

http://archivesinternet.bnf.fr/20170701090810/http://leschibanis.free.fr/Menu1.html



# Parcours guidés

#### Parcours sur un thème

- des sélections réalisées et des textes rédigés par des bibliothécaires et des partenaires institutionnels ou des chercheurs
- Un permalien est associé à chaque source décrite

Travail d'éditorialisation important. Publication en ligne (PDF).

Forme d'entrée dans les collections.



Un parcours guidé sur l'intelligence artificielle



# Application « Archives de l'internet Labs »



#### Recherche plein texte

Recherche simple Recherche avancée Recherche experte Recherche n-gram

### Copie et utilisation des archives

- Possibilité de consulter des vidéos et d'écouter des enregistrements audio
  - prise casque sur les postes publics + casques audio (disponibles en banque de salle)
- Autorisation de copier/coller des extraits de textes à des fins de (courtes) citations
- Toute autre réutilisation des contenus nécessite l'autorisation des ayants droit
- Impossibilité de télécharger des contenus
- Des utilisations plus avancées sont proposées au BnF DataLab

Les services aux chercheurs



# Les dispositifs d'accueil

- Chercheuses et chercheurs associés
   Pour un projet de recherche en lien avec les collections de la BnF
- Accompagnements par le BnF DataLab
   Pour un projet sur les collections numériques (numérisées ou nativement numériques) pouvant avoir une approche expérimentale
  - En répondant à un appel à projet DataLab (une fois par an)
  - En sollicitant le BnF DataLab

Le BnF DataLab travaille en partenariat avec HumaNum

#### Archives du web et recherche - collecte

#### Retrouver des sites disparus

• Approche archivistique « classique » Exemple : site commémoratif centenaire.org

#### S'assurer de l'archivage des contenus étudiés

• Collecte dans le cadre de partenariats Exemple : centre du féminisme d'Angers

#### **Coproduire une collecte**

• Dans le cadre d'un projet de recherche Exemple : collecte sur les littératures francophones



#### Offre de services

- Parcours guidé (médiation)
- Participation aux collectes courantes
- Collecte Corpus de recherche



#### Archives du web et recherche - data

#### Indexer des collections massives

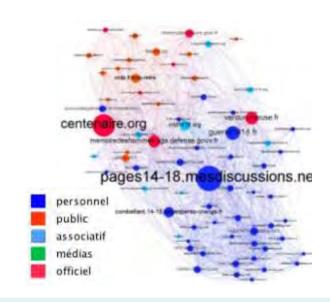
- Indexation d'une collection d'intérêt
- Indexation dans le cadre d'un projet de recherche
- Déterminer quels fichiers doivent être indexés

Exemple: collection Attentats 2015

#### Fouille de données et création de corpus

- Requêtes et extractions à partir des :
- o métadonnées (URL, liens, types MIME)
- index (recherche textuelle)
- mise en place d'un process et accueil d'un IGE/IGR

Exemple : cartographie du web de la Grande Guerre



#### Offre de services

- Indexation
- Extraction de données et métadonnées
- Aide à la fouille de données



#### Portail dédié du BnF DataLab

#### BnF Archives de l'internet Outils pour la recherche

Archives de l'in	ternet Hyp	he BnF	SolrWayback	Vérification d'URL	Sélections	Ressources
Ce portail regroupe un er ne plus fonctionner à tou				lans les archives de l'interne chives de l'internet.	t. Ces outils sont ex	périmentaux et peuvent
H Hyphe BnF	Outil de constitution, d'exploration et de catégorisation de corpus web développé par le médialab de Sciences Po. Hyphe BnF permet de travailler sur les Archives de l'internet.					
SolrWayback	Outil de recherche plein texte, de fouille et de datavisualisation développé en collaboration avec la Bibliothèque Royale du Danemark.					
Q Vérification d'URL	Vérifier la présence d'un site, d'une partie de site, d'une page ou d'un fichier dans les archives à partir de son adresse URL exacte.					
<b>: =</b> Sélections	Rechercher dans les listes de sites sélectionnés pour les collectes ciblées thématiques, régionales et projets.					
Ressources	Supports de formations, jeux de données, documents et liens utiles pour des séances de travail au BnF DataLab.					

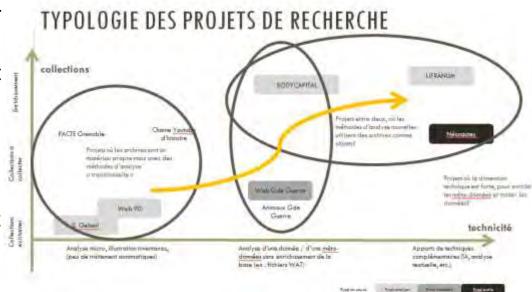
#### Conseils méthodologiques aux chercheurs

#### Les freins à l'usage des archives du web

- Une source difficile à appréhender et à citer
- Une consultation limitée aux emprises de la BnF et des bibliothèques partenaires
- Une faible proportion de collections indexées en plein texte

#### Métadonnées et données dérivées

- ☐ Créer des corpus par étape :
- des métadonnées pour appréhender l'archive
- mixer les approches quantitatives et qualitatives
- créer et exporter des données dérivées
- partager les corpus
- abaisser la marche technique
- tester et expérimenter de nouveaux outils
- faciliter l'approche par la donnée





## Buzz-F Etude de la viralité en ligne



Dispositif d'accueil : projet lauréat AAP DataLab 2021

#### Le projet

- ☐ Etude de la viralité : Analyse des phénomène viraux sur le web et de leur diffusion sur le web français
- ☐ Convention avec l'Université du Luxembourg
- ☐ Equipe C2DH
- Valérie Schafer, historienne
- Fred Pailler, ingénieur d'étude, sociologue
- ☐ Collecte large

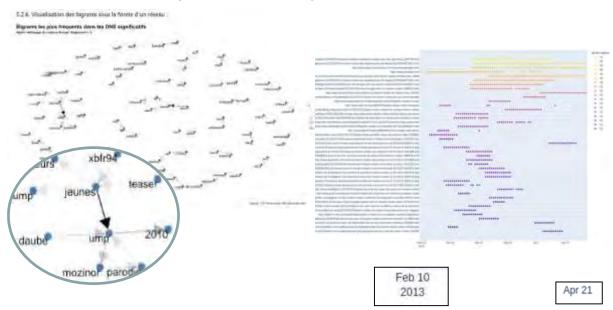
Sites du domaine .fr et hébergés en France. Plusieurs millions de sites crawlés (2000 URL par site). Indexation limitée à l'URL et à la date de capture

☐ Collecte « Actualités » : recherche textuelle



### Buzz-F datavisualisation

### Etude de deux phénomènes : lip dub 2011 et Harlem shake 2013

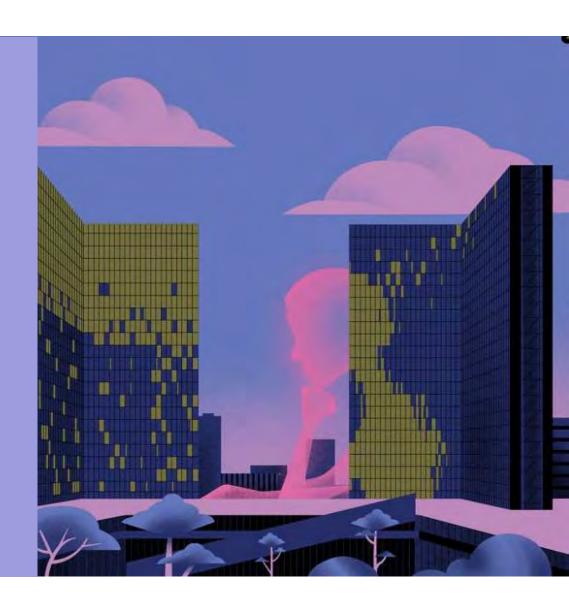


- La viralité : un phénomène difficile à appréhender
- dont l'empreinte se révèle dans les archives du web

- ☐ Une approche des corpus par les données / métadonnées
- ☐ Dépôt et datavisualisation Méthodologie et données sur GitLab Huma-num
- Tokenisation des URL
- Chronologie des traces numériques
- Création de sous-corpus vidéos



L'IA et les archives du web à la BnF



# L'IA et la sélection

- Les correspondants peuvent utiliser des outils IA pour réaliser leur veille
- Une sélection thématique et un parcours guidé sur l'intelligence artificielle

• Mais la sélection reste humaine avec le choix des paramètres de collecte

# {BnF

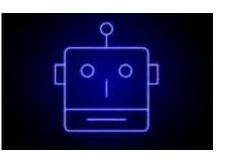


# Impact des nouvelles technologies sur la collecte

Les évolutions technologiques comme les API, l'IA oblige à revoir les processus de collecte.

Par exemple pour des collectes de plateforme, les chercheurs utilisent des API pour récupérer le contenu.

La BnF pratique des formes hybrides comme pour sa collecte Youtube, puisqu'elle récupère les métadonnées de l'API en amont de son crawl. Cela permet de récupérer la liste des vidéos par chaine et d'évaluer le poids de la collecte. Une des conséquences est de devoir recréer une interface complète et artificielle de consultation.



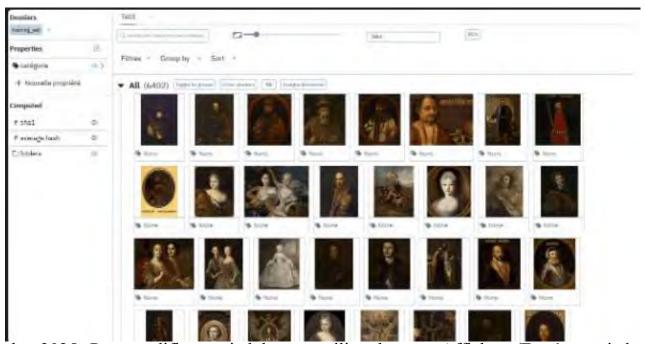
## La collecte à l'heure de l'IA

- Quand des sites embarquent de l'IA, le robot ne fait pas de prompt et pour les sites marchands, le robot n'a pas de préférences de consommation.
- Collecte de vidéos, de témoignages sur l'IA

```
ontal lang="fr-FR" data-build="prod-49b56bec084d41ac16d6a039a3add5463a0e2b9d" dir="ltr" class="">khead>cmarSet="UTF-8"/>kmeta name="viewport" content="width=device-width, initial-scale=1, viewport-fit="".
import " as route@ from "https://cdm.paistatic.com/assets/hbaqc5i44eqo3xfk.js";
import " as route1 from "https://cdm.oaistatic.com/assets/hknra0Ittaxogo0e.js";
import " as route2 from "https://cdn.oaistatic.com/assets/ioxidkm4sd4fiy9u.js";
  window. reactRouterRouteModules = {"root":route0, "routes/ conversation":route1, "routes/ conversation. index":route2};
import("https://cdm.oaistatic.com/assets/byj92clh16ig8xit.js");</script></+$?--><template id="8:2">//template><!--/$--><script nunce="87e4eb65-965f-4a5d-ble9-a40d8693f8fd" id=" R ">requestAnimationFrame(function())[$RT=per
$RC=function(a,b){if(b=document.getElementById(b))(a=document.getElementById(a))?(a.previousSibling.data="$~",$RB.push(a,b),2===$RB.length88("number"!==typeof $RT?requestAnimationFrame($RV.bind(null,$RB)):(a=performance.nc
vscript nonce="87e4eb65-965f-4a5d-ble9-a40d8693f8fd" window.ReactQueryError ??= class ReactQueryError extends Error [];
null==Promise.withResolvers&&(Promise.withResolvers=function(){let e,r;return{promise:new Promise((s,o)=>{e=s,r=o}),resolve:e,reject:r}});
window. REACT QUERY CACHE 23= {}
window. REACT QUERY CACHE ["[\"models\",\"\\\"iBN\\\";false,\\\"isGizmo\\\";false\\"]"] = Promise.withResulvers();
window. REACT QUERY CACHE ["[\"promptStarters\",8,null]"] = Promise.withResolvers();
<div hidden id="5:1" ></div><script nonce="87e4eb65-965f-4a5d-ble9-a48d8693f8fd" >$RC("B:1","5:1")//script>
«script monce"87e4eb65-965f-4a5d-ble9-a40d8893f8fd".window. REACT QUERY CACHE ["[\"promptStarters\",8,null]"].resolve(({items:[[id:"lb3b651a",title:"Créer une image",description:"pour ma présentation",oneliner:"Créer une image description:"pour ma présentation une image description une image 
%/script>
<script nonce="87e4eb65-965f-4a5d-ble9-a40d8693f8fd">window. REACT QUERY CACHE ["[\"models\",\"(\\\"IIM\\\":false,\\\"is6izmo\\\\":false}\\"]"].resolve((h=>({categories:[{color:"#600000",tagline:"",defaultModel:"auto",labe
<div hidden id="5:0"></div><script nance="87e4eb65-965f-4a5d-b1e9-a48d8693f8fd">$RC("B:0"."$:0"\(/script></body></html>
```

# La recherche dans les archives à l'heure de l'IA

- Des perspectives pour les images
  - Panoptic : outil de comparaison d'images, développé par le CERES



# Des perspectives pour les analyses textuelles

- Le projet AdaptMed : mise en place d'un système de génération automatique de reformulation pour les termes médicaux
  - Annotation des contenus collectés pendant le Covid pour constituer une base

# Les limites à la recherche

• Les données des archives du web sont sous droit, ce qui limite les possibilités d'explorer les données avec des outils externes à la BnF

### (BnF Intelligence artificielle

### Feuille de route 2021-2026

À la croisée des chemins, la Bibliothèque nationale de France se constitue à la fois comme un observatoire privilégié, comme une instance réflexive et critique, et comme un moteur des évolutions liées à l'intelligence artificielle. Cette feuille de route présente une projection à 5 ans, fortement ancrée dans les réalisations actuelles et les projets en cours. Davantage qu'un achévement, la démarche proposée est envisagée comme l'initialisation d'une dynamique qui sera appelée à s'étoffer et à évoluer.

#### Actions







Organisar Ia R&D



Acquarir de nouvelles compétences

entispensables à l'émergence et à la conduite des projets (A



Propurer l'infrastructure et les données

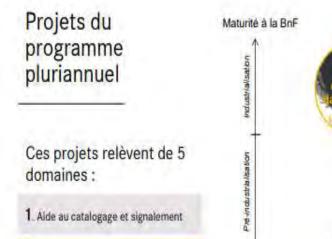


Morter un programme pluritannuel sur FIA en partenarat avec des acteurs clés

#### Jalons



## {BnF



à la conservation

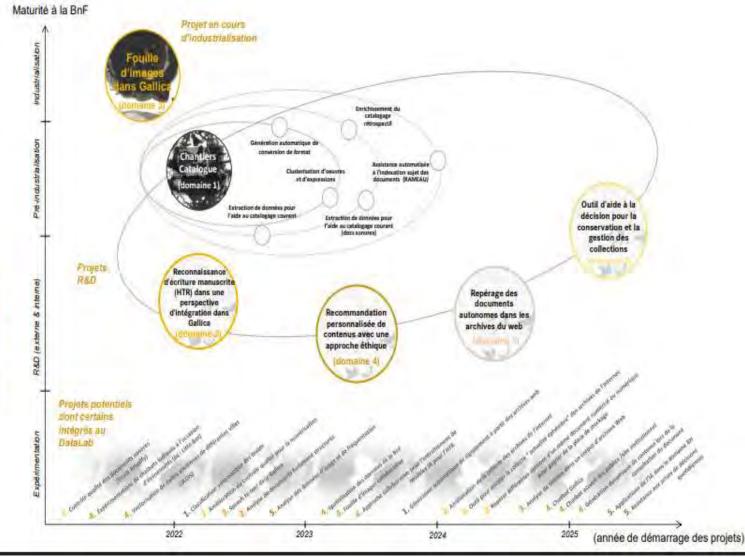
3. Exploration, analyse des collections

2. Gestion des collections, des entrées

- et amélioration de l'accès
- Médiation, valorisation et éditorialisation des collections
- 5. Aide à la décision et au pilotage

# Information et contact www.bnf.fr/fr/feuille-de-route-ia

BIBLIOTHÉQUE NATIONALE DE FRANCE JUIN 2022



# Merci!



• Contacts: depot.legal.web@bnf.fr et datalab@bnf.fr



 Webcorpora et Carnet de recherche de la BnF

https://webcorpora.hypotheses.org/ et https://bnf.hypotheses.org/



• Guide des archives du web

https://www.bnf.fr/sites/default/files/2025-02/Guide Archives internet BnF.pdf



• L'IA à la BnF: https://www.bnf.fr/fr/lintelligence-artificielle-la-bnf