

Recherche d'informations sur Internet (perfectionnement) méthodologie et outils disponibles

TENDANCES 2023



Recherche d'informations sur internet

 Rappels méthodologiques .

Je cherche des pages web, des informations ponctuelles, des personnes, des types de ressources particuliers... .

Je cherche des références bibliographiques .

Je cherche des publications scientifiques (PDF) .

la recherche d'information à l'heure de l'intelligence artificielle .

[support https://framindmap.org/c/maps/446941/public](https://framindmap.org/c/maps/446941/public)

Recherche d'informations sur Internet (perfectionnement) méthodologie et outils disponibles



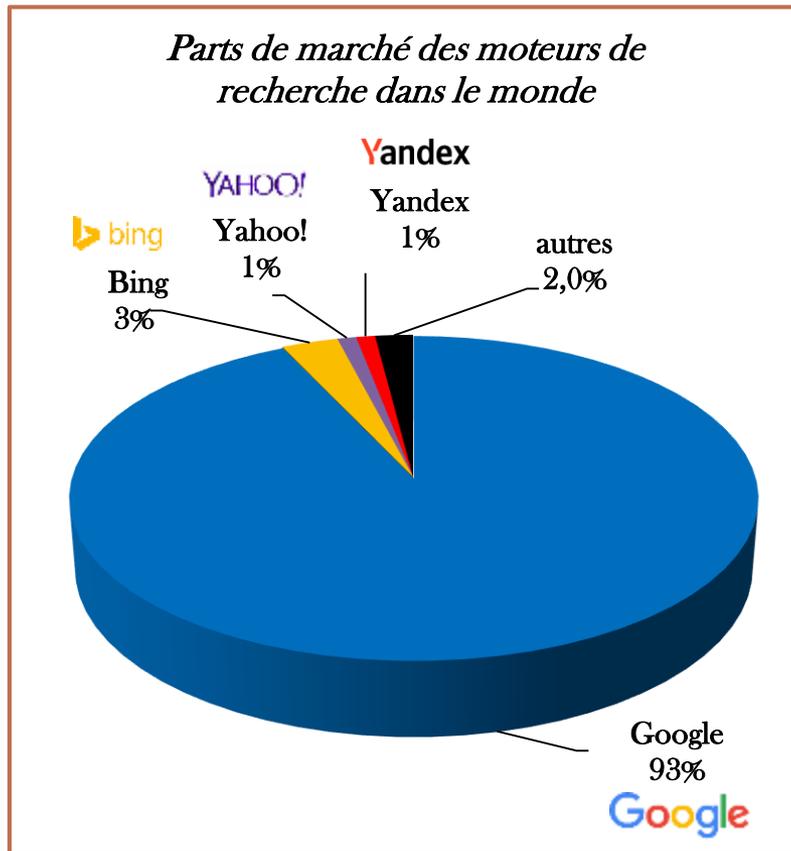
support général

<https://urfist.chartes.psl.eu/ressources/recherche-d-informations-sur-internet-perfectionnement>

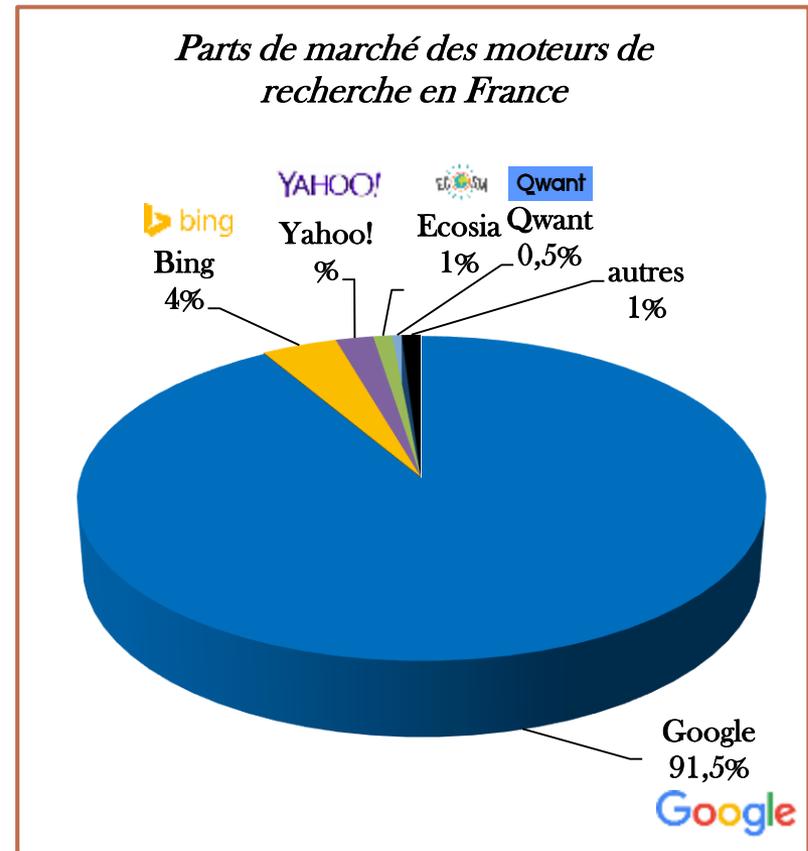
(v. 2021)

Moteurs de recherche

Monde



France



- les moteurs de recherche en 2023

des alternatives aux *Big tech* ?

- Google (et Bing) dominant toujours le marché
- quelques alternatives à tester/utiliser (Brave, DuckDuckGo, Qwant), avec cependant des limites (taille de l'index, couverture des sources francophones, fonctionnalités de recherche inégales...)

modèle économique et données personnelles ?

- des modèles encore souvent dominés par le recours à la publicité et/ou la récupération de données personnelles
- développement d'offres ou de solutions payantes qui doivent néanmoins trouver leur place dans un paysage avant tout gratuit

des fonctionnalités plus poussées

- développement de nouvelles fonctionnalités (ex. recherche sémantique cf. *Knowledge graph* de Google, recherche inversée d'images)
- Google reste cependant dominant dans la panoplie de fonctionnalités avancées (tris, filtres, équations) possibles

des moteurs de recherche aux assistants personnels

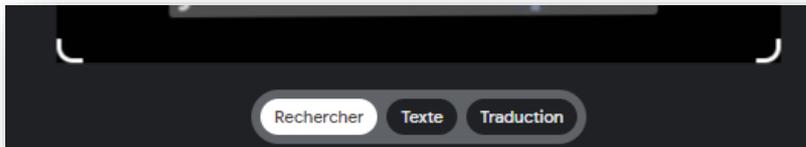
- de moins en moins des moteurs de recherche (de liens) et de plus en plus des moteurs de réponse
- développement d'autres formes d'interactions (ex. : conversation)
- prise en compte de fonctionnalités / contenus issus de l'intelligence artificielle [IA]

Moteurs de recherche

- 2023 : nouveautés Google
mise en valeur et développement de la recherche inversée d'images



intégration directement dans le moteur classique



fonctionnalités avancées :

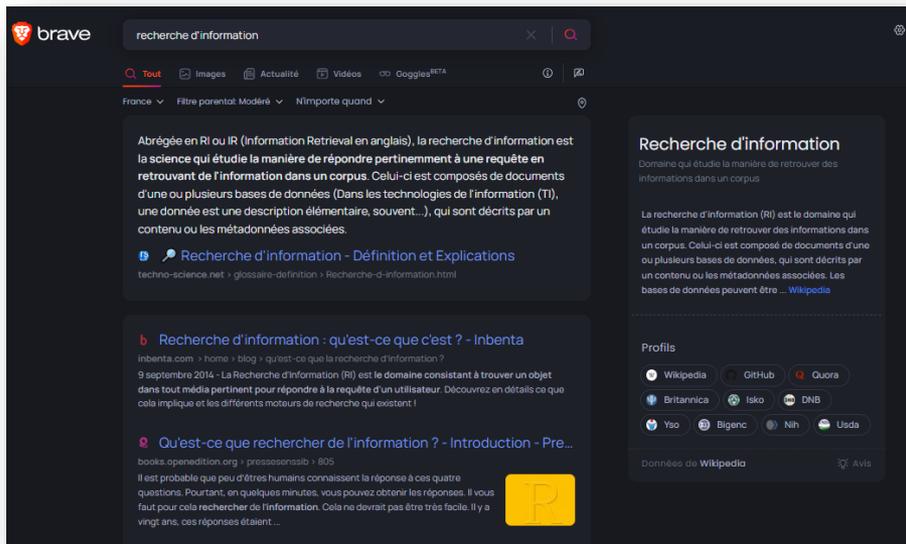
- sélection d'une portion d'image
- sélection de texte sur l'image
- traduction du texte
- remplacement du texte traduit dans l'image d'origine



Moteurs de recherche

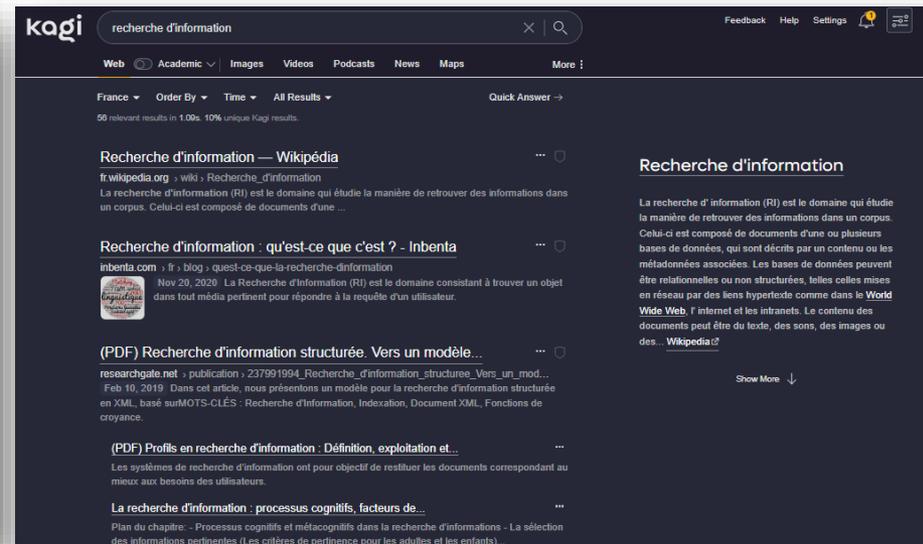
- 2023 : deux *outsiders* à suivre ?

Brave search : <https://search.brave.com/>



- ouvert en 2021, écosystème Brave (navigateur, etc.), fondé par l'un des cofondateurs de Firefox
- index 100 % autonome
mais en cours de constitution et ne proposant pour l'instant ni le filtre vidéo ni le filtre images
- fonction « Goggles » pour personnaliser les sources (version bêta)
- possibilité de supprimer la publicité avec l'offre premium

Kagi : <https://kagi.com/>



- ouvert en 2022, pour un entrepreneur yougoslave
- nécessite un compte (limité à 100 requêtes en version gratuite)
- index : Google et Bing
- très nombreux tris (ex. : par fraîcheur ou site) et filtres (ex. : *academic*, création de *lenses* personnalisées) et possibilité de booster certains types de résultats
- contenus sémantiques / IA : *summarize page* et *ask questions about document*

« Chaque outil dispose d'un corpus spécifique et de fonctionnalités de recherche différentes. [...]

Il est donc crucial d'utiliser **plusieurs moteurs** pour arriver à un résultat satisfaisant.

D'autre part, il ne faut **pas** utiliser tous ces moteurs **de la même façon** »

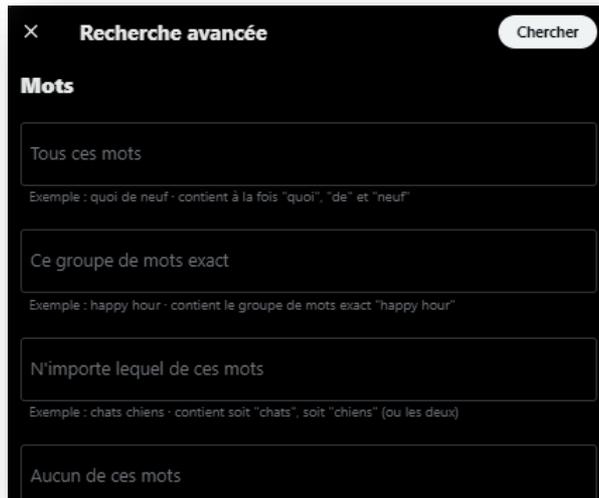
C. Tisserand-Barthole, *Netsources*, n° 149 (11-12/2020)

moteurs de recherche « classiques »	réseaux sociaux
<p data-bbox="426 575 745 611">indexation du web</p> <p data-bbox="423 726 749 762">filtre algorithmique</p> <p data-bbox="396 878 776 913">popularité (<i>PageRank</i>)</p> <p data-bbox="365 1029 807 1065">confiance dans les sources</p>	<p data-bbox="1051 558 1634 644">interrogation du site en temps réel (fraîcheur et tendances)</p> <p data-bbox="1232 726 1450 762">filtre humain</p> <p data-bbox="1108 878 1574 913">sérendipité (signaux faibles)</p> <p data-bbox="1097 1029 1586 1065">confiance dans les personnes</p> <p data-bbox="1012 1243 1673 1329">une indexation souvent incomplète par les moteurs de recherche classiques</p>

Réseaux sociaux

- 2023 : deux réseaux à connaître dans un cadre professionnel et académique

Twitter : <https://twitter.com/>



LinkedIn : <https://www.linkedin.com/>



- le réseau social avec le plus de contenu professionnel ou académique pertinent
- actualité 2022-2023 : rachat par E. Musk
disparition de certaines fonctionnalités
retour / développement de comptes problématiques
mais pas de désertion importante
compte nécessaire pour pouvoir interroger le réseau
- des fonctionnalités de recherche (booléens, filtres, [recherche avancée](#)) et de tris avancées permettant des recherches fines et maîtrisées

- réseau social avec d'importants contenus professionnels ou académiques, mais variable selon les domaines
- compte nécessaire pour pouvoir interroger le réseau
- des fonctionnalités de recherche et de tris peu développées
- phénomène de « boîte noire » (dont résultats non exhaustifs ou aléatoires)
- les abonnements sont pris en compte dans la présentation des résultats

pour aller plus loin sur Twitter : [support de formation](#)

« Quoi qu'on en dise, [Twitter est] toujours à ce jour le réseau social le plus important pour les professionnels de l'information, que ce soit en termes de **contenus** ou de **fonctionnalités**. »

([C. Tisserand-Barthole](#), *Bases*, 01/2023)

- les réseaux sociaux en 2023

des informations complémentaires des moteurs de recherche

- informations en temps réel et actualités, avis et commentaires
- informations personnalisées (abonnements, recherche dans sa communauté)
- fonctionnalités de recherche spécifiques

des services de plus en plus fermés

- nécessité de plus en plus d'avoir un compte pour accéder au contenu et aux outils de recherche
- des algorithmes de plus en plus imposés

le maintien de Twitter et le développement de LinkedIn ?

- Twitter : maintien de la communauté et de la plupart des fonctionnalités de recherche malgré le rachat
- LinkedIn : de plus en plus de contenus professionnels (y compris par des acteurs ayant quitté Twitter)

toujours l'intérêt de l'équation site: sur Google ?

Moteurs de recherche scientifiques

moteurs de recherche « classiques »	moteurs de recherche scientifiques
<p data-bbox="446 554 726 586">indexation du web</p> <p data-bbox="301 753 871 786">fonctionnalités de recherche générales</p> <p data-bbox="295 1008 880 1082">des acteurs généralement commerciaux (Google, Microsoft...)</p>	<p data-bbox="1064 518 1619 596">indexation de corpus spécifiquement académiques</p> <p data-bbox="991 696 1692 861">fonctionnalités et filtres spécifiquement liés aux publications (auteur, titre...), aux types de documents scientifiques (<i>preprints</i>, brevets...) ou aux acteurs (auteurs, institutions...)</p> <p data-bbox="991 1008 1692 1082">des acteurs plus divers (Google, fournisseurs de service, services de documentation)</p>

Page de synthèse : <http://urfist.chartes.psl.eu/ressources/exploiter-l-open-access-en-recherche-d-informations>

Open access et accès libre 2

Je cherche des publications en *open access*..... 13

Le gold open access : revues et ouvrages..... 16

- les revues..... 16
- les plateformes de revues francophones 17
- les ouvrages 18

Le green open access : les plateformes de dépôt et d'autoarchivage..... 18

- les archives ouvertes institutionnelles..... 18
- les archives ouvertes thématiques (exemples les plus connus)..... 19
- le cas des thèses..... 20
- le cas des données de la recherche - *research data* (exemples) 22

Bases bibliographiques et moteurs de recherche scientifiques..... 23

- les bases bibliographiques : Web of Science et Scopus 23
- les moteurs de recherche scientifiques dédiés à l'*open access* (moissonneurs OAI)..... 23
- les moteurs de recherche scientifiques généralistes 25
- les méga-index..... 26
- les extensions de navigateurs 27

De l'open access plus ou moins open et plus ou moins black 32

- demande aux auteurs 32
- les réseaux sociaux académiques 33
- autres réseaux sociaux..... 33
- les archives du web 34
- les bibliothèques clandestines..... 34

Je souhaite publier ou diffuser mes travaux en *open access* 36

Politiques institutionnelles..... 37

- politique des financeurs 37
- positionnement des éditeurs..... 37

Publier en gold open access (revues et ouvrages)..... 38

- publier un article 38
- publier un ouvrage 39
- éditeurs prédateurs et revues pseudo-scientifiques..... 39

Déposer / diffuser en open access (archives ouvertes) 42

- le cas des *preprints* 44
- action individuelle du chercheur..... 45

Documents complémentaires 45

Moteurs de recherche scientifiques

- les moteurs de recherche scientifiques en 2023

des *outsiders* qui cherchent à se démarquer

- un outil gratuit toujours central : Google scholar
- des acteurs plus spécialisés (ex. : *open access* [Brave], différents types de productions scientifiques, ajouts de données bibliométriques [Dimensions, The Lens])

l'importance de diversifier les outils gratuits et payants

- complémentarité des couvertures et des fonctionnalités de recherche
- accès au texte intégral payant
- voir avec son institution

l'intérêt des extensions de navigateurs

l'importance de rester vigilant face à la notion d'*open access*

- quels accès au texte intégral (texte, images)
- quelles versions disponibles (*preprints, postprints...*) pour quelle validation scientifique
- interrogation du texte intégral ou des seules métadonnées

le développement d'outils utilisant de l'intelligence artificielle

- sémantique et clustering
- suggestions
- tâches d'assistance

- évolutions de la recherche internet

recherche multimodale et universelle : élargir les bases de données

- des données de plus en plus nombreuses (*big data*)
- une diversification des contenus

recherche sémantique : se faire comprendre par la machine comme par un humain

- une compréhension des concepts (vocabulaires, entités, relations..., ex. *knowledge graphs*)
- une compréhension du langage naturel (ex. : NLP - *natural language processing*)
- des modèles de langages (ex. : LLM - *large language models*)

recherche contextuelle et conversationnelle : communiquer avec l'humain

- une compréhension du contexte de la requête et de l'intention de l'internaute
- un historique des échanges précédents

+ un entraînement (*deep learning*, réseaux de neurones et travailleurs du clic)

pour comprendre les évolutions en cours



<https://urfist.chartes.psl.eu/ressources/evolutions-des-moteurs-de-recherche-sur-internet>
(2016)

Tendances
de la recherche web
2020



<https://urfist.chartes.psl.eu/ressources/tendances-de-la-recherche-web-2020>
(2020)

Intelligence artificielle

- évolutions des outils de recherche

IA conversationnelle : répondre à partir de bases de données

- chatbots

- agents conversationnels (ex. : Siri, Alexa, Cortana)

IA générative (GenIA) : créer du contenu original à partir d'entraînements pour fournir des contenus similaires aux données d'apprentissage



ChatGPT

Chat <-> chatbot

GPT : *Generative Pre-trained Transformers* (modèles de langage pré-entraînés)

Intelligence artificielle

« ChatGPT est capable de générer des réponses à des questions, de compléter des phrases, de traduire des textes, d'écrire des articles et de tenir des conversations avec des humains.

Il peut également synthétiser des textes suivant un ensemble de contraintes, telles que le ton, le style et le sujet. Il a également été utilisé pour différents usages, tels que la génération de sous-titres pour des vidéos et la création de chatbots. »

[Wikipédia](#)

Intelligence artificielle

moteurs de recherche généralistes	chatbots IA de type ChatGPT
<p>bases de données mises à jour régulièrement</p> <p>approche par mots-clés (et de plus en plus langage naturel)</p> <p>requête (avec booléens, filtres et tris)</p> <p>liste de résultats (et de plus en plus réponses)</p> <p>pas de conservation de l'historique d'une requête à l'autre (itérations)</p> <p>→ moteurs de recherche</p>	<p>données préformatées à partir de bases de données parfois anciennes</p> <p>approche par énoncé (langage naturel)</p> <p>« prompt » plus ou moins précis</p> <p>réponse (parfois sans source)</p> <p>recherche conversationnelle (échanges)</p> <p>→ moteurs de réponses</p>

Intelligence artificielle



Prompt : "photographic dose up portrait of a smiling male robot, hyperrealistic, cinematic lighting" - Negative Prompt : "hair, skin" - Leonardo.ai / DreamShaper v5

l'importance du « prompt »



Prompt : "Portrait of an asiatic man with face tattoos wearing a beanie and earrings studio lighting cinematic angle 85mm 1.4 lit from above ultra detailed" Dall-E

Intelligence artificielle

les « hallucinations » de l'IA



Décris moi le tableau la vierge à la fontaine de Raphaël

 Le tableau "La Vierge à la fontaine" de Raphaël est un tableau célèbre de la Renaissance italienne. Il représente la Vierge Marie et l'enfant Jésus assis sur un banc de pierre près d'une fontaine, dans un paysage bucolique. La Vierge tient l'enfant Jésus sur ses genoux et lui tend un verre d'eau de la fontaine, tandis que l'enfant boit avidement.

Le tableau est remarquable pour sa composition soignée et sa luminosité. La Vierge et l'enfant sont placés au centre du tableau, entourés d'un paysage verdoyant et de nuages dans le ciel. La Vierge est vêtue de vêtements luxueux et porte une robe bleue et une cape rouge. L'enfant Jésus, avec son visage souriant et son corps robuste, est représenté comme un enfant en bonne santé et en pleine croissance.

"La Vierge à la fontaine" est actuellement conservé au musée du Louvre à Paris. C'est une œuvre importante de Raphaël, un peintre et architecte italien de la Renaissance, connu pour ses tableaux de la Vierge Marie et de l'enfant Jésus.

- outils de recherche IA en 2023

de nouvelles manières de chercher de l'information

- recours au langage naturel et à l'incrémentation des résultats précédents
- une perte de maîtrise (absence de tris / filtres explicites)

des tentatives généralement prudentes

- des fonctionnalités non stabilisées
- de l'IA sur certains points seulement

beaucoup trop d'informations incorrectes voire fausses

- phénomène de boîte noire (sources utilisées ?, algorithmes de classement ?, exhaustivité ?)
- des réponses non sourcées
- des réponses et des citations biaisées voire inventées

des réponses généralement superficielles

- de simples copier-coller ou paraphrases de contenus internet
- une prépondérance des premiers résultats des moteurs de recherche/bases de données bibliographiques

→ ne peuvent suffire seuls pour un état de l'art ou une revue systématique

→ nécessitent de savoir se poser des limites

Intelligence artificielle

- point de vue du documentaliste

« L'IA générative oui, mais pas pour trouver l'information dans les moteurs. [...] L'IA générative a un rôle à jouer dans la recherche d'information, mais cela ne pourra être qu'en amont en phase de recherche de mots-clés, par exemple ou après avoir identifié les informations et les sources pertinentes. »
([C. Tisserand-Barthole, Bases, 04/2023](#))

- point de vue du chercheur

« Although ChatGPT provides succinct answers to most queries, its lack of reliability, transparency, and up-to-date knowledge compared with conventional search engines is a major drawback, particularly for health-related queries.» ([F. Eggmann et M. B. Blatz, Compendium..., 04/2023](#))

- la recherche d'information en 2023
 - un outil toujours central : Google
 - une diversification des autres acteurs et des fonctionnalités qui évoluent (moteurs de recherche « classiques » + réseaux sociaux + moteurs de recherche scientifiques + outils de recherche IA)
 - des questions sur les modèles économiques et la pérennité des services
 - des questions sur les résultats (algorithmiques, chronologiques, personnalisés)
 - de nouvelles manières de rechercher (langage naturel)
 - une interrogation sur la plus-value de l'intelligence artificielle
- diversifier les outils
- être de plus en plus vigilant sur les sources couvertes et les résultats (crédibilité, classement, fraîcheur, pertinence, exhaustivité)

Pour aller plus loin

- méthodologie et outils
 - support de la formation : Aline Bouchard. Recherche d'informations sur internet (perfectionnement). MAJ 2021, à compléter par la carte : <https://framindmap.org/c/maps/446941/public>.
<https://urfist.chartes.psl.eu/ressources/recherche-d-informations-sur-internet-perfectionnement>.
 - Véronique Mesguich. *Rechercher l'information stratégique sur le web. Sourcing, veille et analyse à l'heure de la révolution numérique*. 2^e éd. ADBS-éd. De Boeck, 2021. 236 p. (« Information et stratégie »).
- pour suivre l'actualité du domaine, les incontournables :
 - ressources de FLA Consultants, plus ou moins gratuites : revues *Bases* et *Netsources*, brèves : <https://www.bases-netsources.com/> et compte Twitter : <https://twitter.com/BasesNetsources>
 - blog d'Aaron Tay, <https://musingsaboutlibrarianship.blogspot.com>