

Wikipédia, quinze ans de recherches

LE MONDE SCIENCE ET TECHNO | 11.01.2016 à 16h39 • Mis à jour le 15.01.2016 à 10h56 | Par David Larousserie
(/journaliste/david-larousserie/)



Mario Wagner

Quel succès! Quinze ans après son lancement, le 15 janvier 2001, par les Américains Jimmy Wales et Larry Sanger, l'encyclopédie en ligne Wikipédia reste le premier site non commercial du Web mondial, toujours dans le top 10 des sites les plus fréquentés avec près de 500 millions de visiteurs uniques par mois pour plus de 250 éditions linguistiques. 36,9 millions d'articles sont rédigés, corrigés, améliorés par quelque 2 millions de contributeurs. 800 nouvelles entrées en anglais sont ajoutées chaque jour, 300 en français. La version française tenant la troisième position, avec plus de 1,7 million d'articles, derrière l'anglophone (plus de 5 millions) et la germanique (1,8 million).

Mais Wikipédia, c'est moins connu, est bien plus qu'une encyclopédie qu'on consulte pour se documenter ou faire ses devoirs scolaires. Elle est devenue aussi un objet de recherche en tant que tel, à l'instar d'une tribu d'Amazonie, d'un programme informatique ou d'un patient. La base de données Scopus, l'une des trois plus importantes du monde, recense ainsi plus de 5400 articles ayant pour sujet ou pour objet Wikipédia publiés dans des revues, des actes de colloques ou des livres. Quatorze brevets mentionnent même le célèbre site, selon la même Scopus.

Lire aussi [Rencontre avec les petites mains anonymes qui font Wikipédia](#)

(/pixels/article/2016/01/15/rencontre-avec-les-petites-mains-anonymes-qui-font-wikipedia_4847756_4408996.html)

Autre preuve de l'intérêt académique pour le sujet, en juin 2013, à Paris, se tenait un colloque, coorganisé par le CNRS et le CNAM et intitulé « Wikipédia, objet scientifique non identifié », avec sociologues, spécialistes de sciences de la communication, informaticiens...

Depuis 2011, la Fondation Wikimedia, qui héberge Wikipédia, s'est même

LA FONDATION
WIKIMÉDIA, QUI
HÉBERGE
WIKIPÉDIA, S'EST
DOTÉE D'UN
GROUPE DE
RECHERCHE

dotée d'un groupe de recherche. « *A l'origine, nous étions là pour fournir des outils d'analyse à la communauté. Maintenant, nous sommes un vrai département de recherche, testant de nouvelles technologies et collaborant avec les universités* », résume Dario Taraborelli à la tête des dix personnes de ce département, en Californie. Fin novembre 2015, il annonçait ainsi un projet d'intelligence artificielle capable d'estimer la probabilité qu'une modification soit dommageable ou non à un article et donc susceptible d'être retirée. Auparavant, les systèmes automatiques de détection avaient tendance à trop souvent écarter les contributions pourtant bienveillantes,

freinant l'entrée de nouveaux contributeurs. Le groupe essaie aussi de réduire les asymétries de contenus entre différentes langues en proposant automatiquement des articles à rédiger aux contributeurs des langues minoritaires.

Mais que font tous les autres chercheurs en tripatouillant Wikipédia? De récentes publications témoignent du large spectre couvert. Depuis novembre, une équipe japonaise s'est servie des articles de l'encyclopédie pour analyser les suicides de personnalités dans son pays. Des Britanniques ont construit automatiquement un glossaire technique. Des Turcs ont utilisé le site pour repérer à grande échelle des entités dans des corpus de leur langue. Des Français ont proposé un classement des universités reposant sur les citations des établissements au sein de plusieurs versions linguistiques de Wikipédia. Citons encore un article paru en mai, qui prévoit les pics d'apparition de la grippe grâce aux statistiques de visites des pages de l'encyclopédie.

1,7 million de pages d'articles

Les raisons d'un tel engouement sont simples à comprendre. L'objet est vaste, une quinzaine de gigaoctets de textes (pour la version anglaise). D'utilisation gratuite, contrairement aux données de Facebook, Google ou Twitter, pourtant gigantesques et fournies gracieusement par leurs utilisateurs. Même les données de fréquentation sont disponibles pour chaque article! Les archives sur quinze ans permettent d'avoir du recul historique, tout en ayant un objet toujours rafraîchi. Des versions en plus de 200 langues ouvrent des perspectives pour des comparaisons ou des analyses culturelles. L'ouverture et la transparence offrent aussi ce que les chercheurs adorent : la vérifiabilité et la reproductibilité. Pour parfaire leur bonheur, l'encyclopédie, tel un iceberg, recèle plus de trésors que sa seule vitrine d'articles. Si la version française contient 1,7 million de pages d'articles, elle contient 4,5 fois plus de pages pour les historiques, les discussions et autres coulisses qui font le dynamisme et la réputation du site.

UNE SORTE DE
BAC À SABLE
DANS LEQUEL
S'ÉBROUENT LES
SPÉCIALISTES DU
TRAITEMENT
AUTOMATIQUE DU
LANGAGE

Du coup, presque tous les domaines sont couverts. La sociologie, bien sûr, fascinée par cette démocratie d'un nouveau genre, car auto-organisée et reposant sur quelques règles et le consensus. Les chercheurs, profitant de la transparence du site, y ont également étudié le rôle des « vandales » et autres « trolls » qui mettent leurs pattes malveillantes dans les articles. Les inégalités hommes-femmes particulièrement criantes, avec moins de 10 % de contributrices à l'encyclopédie, ont également donné lieu à beaucoup de littérature et de controverses.

Wikipédia est devenu une sorte de bac à sable dans lequel s'ébrouent les spécialistes du traitement automatique du langage qui disposent là d'un corpus immense pour tester leurs logiciels de reconnaissance de texte, de traduction, d'extraction de sens... C'est aussi le jouet de physiciens, statisticiens, informaticiens... prompts à dégainer leurs outils d'analyse pour en extraire de nouvelles informations ou aider à les visualiser.

Lire aussi [Les coulisses de l'encyclopédie en ligne](#) (/sciences/article/2016/01/11/les-coulisses-de-l-encyclopedie-en-ligne_4845298_1650684.html)

« *Après quinze ans, l'intérêt des chercheurs est toujours là. La première phase était très active car l'objet était nouveau. Cela a contribué à l'émergence de nouveaux domaines comme la sociologie quantitative ou l'informatique sociale*, rappelle Dario Taraborelli. Puis, à partir de 2007, l'apparition

de nouveaux *médias* sociaux a détourné un peu les recherches, avant un renouveau depuis 2010. Notamment parce que nous sommes le seul site important à *publier* nos données quotidiennes de trafic. »

« Notre ambition est de rendre encore plus intelligents les ordinateurs »

Ce renouveau est aussi tiré par une révolution à venir. Wikipédia est devenu l'un des maillons indispensables à un projet particulièrement ambitieux : rassembler toute la connaissance mondiale et la rendre intelligible par des machines. « Notre ambition est de rendre encore plus intelligents les ordinateurs afin qu'ils soient toujours plus utiles à l'humanité », s'enthousiasme Fabian Suchanek, enseignant-chercheur à Télécom ParisTech et artisan de cette évolution qui vise à transformer Wikipédia et d'autres riches corpus en une source accessible aux ordinateurs.

DERRIÈRE DES
PROUESSES QUI
N'ONT L'AIR DE
RIEN SE CACHENT
DE NOUVEAUX
OBJETS : LES
BASES DE
CONNAISSANCE

De tels changements sont en fait déjà à l'œuvre, discrètement. Dans les *moteurs de recherche* par exemple, lorsque l'utilisateur tape un nom de célébrité, apparaissent toujours une liste de liens mais aussi un encadré résumant la biographie de la personne cherchée. Et cela automatiquement : le programme a compris où, dans la page Wikipédia, se trouve l'information souhaitée. Mieux. On peut désormais *poser* des questions explicites, en langage naturel, à ces moteurs : quand Elvis Presley est-il mort ? Où ? Quel est l'âge de François Hollande ?... et *recevoir* des réponses directes, sans *avoir à lire* la page contenant l'information.

Derrière ces prouesses qui n'ont l'air de rien se cachent de nouveaux objets : les bases de connaissance. Les plus célèbres sont Yago, DBpedia, Freebase ou Wikidata. Toutes se sont construites en triturant Wikipédia. Et, preuve des enjeux économiques, les plus grands du Web actuel investissent dans ces constructions. En 2010, Google a ainsi racheté Freebase, qui lui sert pour son Knowledge Graph, l'encadré qui fournit des réponses directes aux requêtes. L'entreprise soutient également financièrement Wikidata, une initiative de la fondation Wikimédia. Amazon a racheté EVI en 2012, anciennement connue sous le nom de True Knowledge, une base de connaissances.

Lire aussi Wikipédia, « un objet scientifique non identifié » (/sciences/article/2016/01/11/wikipedia-un-objet-scientifique-non-identifie_4845331_1650684.html)

En outre, derrière les assistants personnels vocaux des mobiles, Siri, Cortana ou Google Now, se cachent aussi ces fameuses bases de connaissances. Pour *gagner* au jeu Jeopardy en 2011, l'ordinateur Watson d'IBM a bien sûr assimilé bon nombre de données, en particulier de Wikipédia, mais dans une forme prédigérée fournie par la base de connaissances Yago.

Le sujet de ces bases ou graphes de connaissances est très actif. Le chercheur le plus prolifique sur Wikipédia, toutes activités confondues selon Scopus, est par exemple l'Allemand Gerhard Weikum de l'Institut Max-Planck de Sarrebruck, à l'origine de la première base de connaissances, Yago, en 2007. Le second est un Hollandais, Maarten de Rijke, professeur d'informatique à l'université d'Amsterdam, dont les récents travaux utilisent ces graphes. Il est capable de *savoir* de quoi parle un tweet en repérant les noms et les faits à l'intérieur et en les confrontant à Yago ou DBpedia. Il enrichit aussi les émissions de télévision automatiquement en fournissant des liens sur les tablettes ou téléphones, choisis en fonction du thème de l'émission, déterminé grâce aux bases de connaissances.

« ON PEUT FAIRE
DES CHOSES QUI
ÉTAIENT
IMPOSSIBLES
AUPARAVANT »

« Avec ces bases de connaissances, on peut faire des choses qui étaient impossibles auparavant », estime Fabian Suchanek, cofondateur de Yago. Par exemple ? « Extraire de l'information du quotidien Le Monde : combien de femmes en *politique* au cours du temps ? Quel est l'âge moyen des politiciens ou des chanteurs cités ? Quelles compagnies étrangères sont mentionnées ? », énumère ce chercheur en citant un *travail* publié en 2013 avec la collaboration du journal. Le *New York Times* construit sa propre base de connaissances tirées des informations de ses articles. Autre

exemple, il devient possible de poser des questions aussi complexes que : qui sont les politiciens également scientifiques nés près de Paris depuis 1900 ? Ou, plus simplement, quelle est la part des femmes scientifiques dans Wikipédia ?

Mais quelle différence entre ces objets et une base de données ou même une page Wikipédia ? Si un humain comprend que dans la phrase « Elvis Presley est un chanteur né le 8 janvier 1935 à Tupelo, Mississippi », il y a plusieurs informations sur son métier, sa date et son lieu de naissance, une machine ne le comprend pas, et ne peut donc **répondre** à la question simple, pour un humain, « Quand Elvis est-il né ? ». « *C'est un peu paradoxal, mais pour un informaticien, notre langage n'est pas structuré et donc un ordinateur ne peut le comprendre !* », souligne ironiquement Fabian Suchanek. Il faut donc transformer les pages en les structurant différemment, en commençant par repérer les entités, les faits et les relations entre eux. Presley est une entité. Sa date de naissance ou son métier sont des faits. « Né le » et « a pour métier » sont les relations. Tout cela peut **être** codifié en langage informatique.

Une autre particularité de ces objets est qu'ils ne répertorient pas ces faits et entités dans des tableaux, comme la plupart des bases de données, mais en les organisant en arborescences ou en graphes. Les branches correspondent aux liens entre les entités et les faits. Les informaticiens et mathématiciens ont bien sûr développé les techniques pour **interroger** ces graphes et y faire des calculs comme dans un vulgaire tableur. Aujourd'hui, Yago « sait » plus de 120 millions de choses sur 10 millions d'entités (**personnalités**, organisations, **villes** ...).

Réseau reliant des faits et des entités

PETIT À PETIT SE
TISSE UN RÉSEAU
RELIANT DES
FAITS ET DES
ENTITÉS

L'avantage-clé est que le rapprochement devient plus simple entre plusieurs bases de connaissances, celles construites sur Wikipédia mais aussi d'autres concernant les musiciens, les coordonnées GPS, les gènes, les auteurs... Le site [Linkeddata.org](http://linkeddata.org) recense ces nouvelles bases et les liens entre elles. Petit à petit se tisse un réseau reliant des faits et des entités, alors que, jusqu'à présent, la Toile connecte des pages ou des documents entre eux. Cela contribue au rêve de ce que Tim Berners-Lee, le physicien à l'origine du Web, a baptisé « Web sémantique » en 2001. « *Les défis ne manquent pas. La troisième version de Yago est sortie en mars 2015. Nous avons déjà traité la question du temps. Nous traitons aussi plusieurs langues. Il faut maintenant s'attaquer aux "faits mous", c'est-à-dire moins évidents que les dates et lieux de naissance, les métiers, le genre...*, estime Fabian Suchanek. *En outre, tout ne peut pas se mettre dans un graphe !* »

Bien entendu, faire **reposer** la connaissance future de l'humanité sur Wikipédia n'a de sens que si ce premier maillon est solide. La crédibilité de l'encyclopédie a donc été parmi les premiers sujets d'études. Dès 2005, *Nature* publiait un comparatif entre l'encyclopédie en ligne et sa « concurrente » *Britannica*, qui ne montrait pas d'énormes défauts pour la première. D'autres études ont été conduites depuis pour **estimer** l'exactitude, en **médecine** par exemple, Wikipédia étant l'un des premiers sites consultés sur ces questions. Les résultats sont bien souvent satisfaisants.

« *C'est finalement un peu une question vaine scientifiquement, car les comparaisons sont souvent impossibles. On confronte les articles tantôt à des encyclopédies, tantôt à des articles de revues scientifiques...* », estime Gilles Sahut, professeur à l'École supérieure du professorat et de l'éducation, de l'université **Toulouse** -Jean-Jaurès. « *La question a un peu changé de nature. Il faut passer d'une appréciation globale à une appréciation au cas par cas, et donc éduquer afin d'être capable de dire si un article semble biaisé ou complet* », précise ce chercheur, qui a soutenu une thèse en novembre 2015 sur la crédibilité de Wikipédia. Il adosse ce constat à une étude menée sur plus de 800 jeunes entre 11 et 25 ans, pour tester la confiance accordée à l'encyclopédie. Celle-ci s'érode avec l'âge et le niveau de scolarité, mais elle remonte dès lors que les élèves participent. « *Ils découvrent d'ailleurs, comme leur enseignant, qu'il n'est pas si facile d'écrire dans Wikipédia !* », sourit le chercheur en faisant allusion aux difficultés à **entrer** dans la communauté. « *Certes les wikipédiens sont des maîtres ignorants sur les savoirs, comme le dit le sociologue Dominique Cardon, mais ils sont très savants sur les règles et les procédures !* »