

Qu'est-ce que Common Crawl ? Une histoire de l'ensemble de données du web ouvert

By rankstudio.net Publié le 27 octobre 2025 48 min de lecture



Résumé

Common Crawl est une fondation à but non lucratif 501(c)(3) (fondée en 2007) qui maintient un dépôt gratuit et ouvert de données d'exploration web (Source: commoncrawl.org) (Source: commoncrawl.org). Sa mission est de démocratiser l'accès à l'information web en fournissant gratuitement des jeux de données d'exploration web à l'échelle du pétaoctet. Au cours des 15 dernières années et plus, Common Crawl a collecté de l'ordre de 300 à 400 milliards de pages web, couvrant plus de 15 ans d'exploration continue (Source: commoncrawl.org) (Source: www.96layers.ai). Chaque mois, elle ajoute environ 3 à 5 milliards de nouvelles pages (environ 90 To compressés, ~400 To non compressés) (Source: www.96layers.ai) (Source: commoncrawl.org). Bien qu'il ait commencé comme un petit projet (seulement quelques employés) (Source: www.96layers.ai), le corpus publiquement disponible de Common Crawl sous-tend désormais un large éventail d'utilisations de recherche et commerciales. Notamment, il fournit la majeure partie des données d'entraînement pour les grands modèles linguistiques (LLM) modernes – par exemple, plus de 80 % des tokens du GPT-3 d'OpenAl provenaient des données de Common Crawl (Source: www.mozillafoundation.org) – et est cité dans plus de 10 000 publications universitaires (Source: commoncrawl.org) (Source: dallascard.github.io). Il a permis à des startups (par exemple, TinEye, Lucky Oyster) et à des projets de recherche (par exemple, les word embeddings GloVe, l'analyse de la censure web) qui, autrement, n'auraient pas eu les ressources nécessaires pour explorer l'ensemble du web. Common Crawl sert ainsi d'« infrastructure neutre et à but non lucratif » pour les données web (Source: www.96layers.ai), uniformisant les règles du jeu afin que même les petites organisations et les chercheurs puissent accéder à des informations à l'échelle du web.

Ce rapport fournit un **historique et une analyse complets de Common Crawl**. Il couvre les origines du projet (motivations clés, parcours du fondateur, développement précoce), la structure organisationnelle et le financement, les méthodes et la technologie de collecte de données, la croissance du jeu de données, et les **multiples façons dont les données sont utilisées aujourd'hui** (dans l'entraînement d'IA/LLM, la recherche universitaire, les produits industriels, etc.). Nous examinerons le contexte social et technique (par exemple, la domination de Google et le besoin de données web ouvertes), résumerons les **statistiques**



quantitatives (pages collectées, volume de données, nombre de citations) et présenterons des études de cas illustrant l'impact de Common Crawl. Nous discuterons également des défis (biais de couverture, problèmes de droits d'auteur) et des orientations futures. Toutes les affirmations et tous les faits sont étayés par des sources faisant autorité de l'organisation Common Crawl, des médias, des entretiens et des publications de recherche.

Introduction et Contexte

Le **World Wide Web** est devenu un vaste écosystème d'information décentralisé. Les <u>moteurs de recherche modernes comme Google et Bing</u> explorent continuellement le web pour créer leurs propres index, mais ces index sont propriétaires. Au milieu des années 2000, il n'existait **aucun dépôt majeur de données d'exploration web accessible au public** pour les acteurs extérieurs. Seules quelques organisations — notamment l'organisation à but non lucratif <u>Internet Archive</u> — tentaient de préserver les pages web (par exemple, via la Wayback Machine). Cependant, la *Wayback Machine* d'Internet Archive est conçue pour l'archivage instantané à la demande et la navigation de pages web au fil du temps ; elle n'est pas optimisée pour l'analyse de données à grande échelle ou l'exploration algorithmique du contenu du web (Source: <u>dallascard.github.io</u>).

Dans ce contexte, l'idée de construire un « **index web ouvert** » a commencé à émerger. Entrepreneurs et chercheurs ont reconnu que seules les plus grandes entreprises (Google, Microsoft, Yahoo, Baidu, etc.) avaient les ressources nécessaires pour explorer des milliards de pages à haute fréquence, laissant les plus petits acteurs sans accès à ces données brutes. Par exemple, les chercheurs universitaires et les startups avaient souvent besoin de grands corpus web pour le traitement du langage naturel (TLN), l'exploration de données et les tâches d'apprentissage automatique, mais manquaient des moyens d'explorer l'ensemble du web eux-mêmes. Un dépôt ouvert de données d'exploration web **démocratiserait l'accès** et favoriserait l'innovation, de la même manière que les jeux de données ouverts (par exemple, Wikipédia) ont alimenté de nouvelles recherches.

Common Crawl a été conçu et lancé pour répondre à ce besoin. Son fondateur, **Gil Elbaz**, est un entrepreneur en série et technologue : à la fin des années 1990, il a cofondé Applied Semantics (l'entreprise qui a construit la technologie connue plus tard sous le nom de Google AdSense) (Source: www.96layers.ai) (Source: www.96layers.ai). Après l'acquisition d'Applied Semantics par Google, Elbaz a travaillé chez Google jusqu'en 2007. Dans des entretiens, il a expliqué que son départ était motivé par l'inquiétude concernant la concentration des données et son impact sur l'innovation. Il considérait l'exploration propriétaire massive de Google comme la clé de son monopole sur l'innovation en matière de recherche (Source: www.96layers.ai) (Source: www.96layers.ai). Pour contrebalancer cela, Elbaz a imaginé des « sociétés de données neutres » — des projets d'infrastructure ouverts et à but non lucratif qui exploreraient le web et fourniraient les données **gratuitement** aux chercheurs et aux entreprises. L'un de ces projets était **Common Crawl.** fondé en 2007. Comme l'a dit Elbaz :

« Common Crawl était censé être comme une infrastructure neutre et à but non lucratif qui devrait imiter la façon dont Google explorait le web... et ensuite rendre ces données disponibles à tous gratuitement, afin d'uniformiser les règles du jeu du développement technologique » (Source: www.96layers.ai).

La motivation d'Elbaz était donc explicitement d'uniformiser les règles du jeu. Il voulait que les petites startups et les chercheurs universitaires aient les mêmes informations brutes d'« index de recherche » que Google – afin que l'innovation ne soit pas monopolisée par une seule entreprise (Source: www.novaspivack.com) (Source: www.96layers.ai). Cette vision a trouvé un écho auprès d'autres leaders de la communauté du web ouvert. Des technologues éminents tels que Nova Spivack (un des premiers entrepreneurs de l'Internet) et Carl Malamud (un pionnier des données gouvernementales ouvertes) ont rejoint le conseil d'administration fondateur de Common Crawl (Source: www.novaspivack.com). Au fil du temps, le conseil consultatif s'est élargi pour inclure des personnalités éminentes comme le directeur de recherche de Google, Peter Norvig, et le directeur du MIT Media Lab, Joi Ito (Source: www.thekurzweillibrary.com) (Source: commoncrawl.org), soulignant l'importance du projet.

En quelques années, Common Crawl est devenue une fondation indépendante à but non lucratif. Dès son lancement, elle a été enregistrée en tant qu'organisation 501(c)(3) californienne, la **Common Crawl Foundation** (Source: commoncrawl.org) (Source: commoncrawl.org). Son énoncé de mission est succinct : « démocratiser l'accès à l'information web en produisant et en maintenant une exploration ouverte du web ». La page d'accueil de Common Crawl la décrit comme « un dépôt gratuit et ouvert de données d'exploration web qui peut être utilisé par n'importe qui » (Source: commoncrawl.org). Gil Elbaz a occupé le poste de Président du Conseil d'administration et est souvent reconnu comme le fondateur du projet (Source: commoncrawl.org) (Source: www.novaspivack.com). Parmi les autres membres clés de l'équipe initiale figuraient l'ingénieur principal **Ahad Rana** et plus tard la directrice **Lisa Green** (anciennement de Creative Commons) (Source: www.novaspivack.com).



Structure Organisationnelle et Financement

Common Crawl fonctionne comme une petite organisation à but non lucratif. Sa page d'accueil et ses pages d'équipe de 2025 indiquent que l'équipe principale a toujours été très petite — littéralement « moins de cinq personnes » dans les premières années (Source: www.96layers.ai). Par exemple, au début des années 2010, le projet fonctionnait avec seulement une poignée d'ingénieurs et de bénévoles. Même au moment où OpenAl a publié le document GPT-3 en 2020, Common Crawl n'aurait eu qu'un seul employé à temps plein (Source: www.96layers.ai) (bien qu'en 2025 l'équipe soit plus grande). Gil Elbaz est Président (et a été co-Président de Factual/Foursquare), et des noms comme Peter Norvig sont conseillers (Source: commoncrawl.org). Cependant, les opérations quotidiennes reposent sur une très petite équipe permanente et les contributions de bénévoles et de collaborateurs.

L'organisation est financée principalement par des **dons et des parrainages**, en particulier des fournisseurs de services cloud. À partir de 2012, Amazon Web Services (AWS) a hébergé les données de Common Crawl sans frais dans le cadre du programme AWS Public Datasets (Source: <u>alchetron.com</u>). Le parrainage de données publiques d'AWS fournit l'immense stockage requis (plusieurs centaines de téraoctets) sans facturer Common Crawl. D'autres plateformes cloud (par exemple, Microsoft Azure, Google Cloud) peuvent également être impliquées dans les archives, mais AWS est l'hôte principal. De plus, des entreprises comme Amazon ont offert des concours de petites subventions (par exemple, 50 \$ de crédits AWS) pour encourager l'utilisation des données (Source: <u>commoncrawl.org</u>). La fondation reçoit probablement aussi de modestes dons philanthropiques, bien que **Common Crawl n'ait jamais reçu d'investissement en capital-risque ni fonctionné comme une entreprise commerciale**. (Elle reste délibérément une organisation à but non lucratif pour rester « neutre » et exempte de motifs de profit (Source: <u>www.novaspivack.com</u>) (Source: <u>www.96layers.ai</u>).)

En bref, Common Crawl est le produit collaboratif de quelques technologues passionnés et de l'écosystème du cloud computing. Ses coûts d'exploitation relativement bas (car elle contourne les frais de stockage) lui permettent de persister avec un financement minimal. En 2024, Common Crawl reste « largement inconnue du grand public », mais elle est reconnue pour jouer « un rôle important » dans des domaines comme l'IA générative (Source: www.mozillafoundation.org). Le rapport 2024 de la Mozilla Foundation souligne que Common Crawl est « une petite organisation à but non lucratif » avec un impact massif (Source: www.mozillafoundation.org).

Collecte de Données : Exploration et Technologie

Common Crawl exécute un robot d'exploration web automatisé (nommé **CCBot**) qui scanne continuellement le web public pour construire son jeu de données. Le robot d'exploration est construit sur le cadre open source <u>Apache Nutch</u>, qui gère la découverte d'URL, la récupération de pages et le suivi des hyperliens (Source: <u>datadome.co</u>). (En fait, en 2013, Common Crawl est passé à l'utilisation d'Apache Nutch comme robot d'exploration principal « au lieu d'un robot d'exploration personnalisé » (Source: <u>alchetron.com</u>), et il a migré de l'ancien format de fichier « ARC » vers le format standard **WARC** en même temps (Source: <u>alchetron.com</u>).) CCBot s'identifie dans l'agent utilisateur comme « CCBot/2.0 » (Source: <u>datadome.co</u>), bien qu'il soit déconseillé de se fier uniquement à la chaîne d'agent utilisateur car les robots peuvent usurper des identités. CCBot explore à partir d'adresses IP Amazon AWS. Au cours des premières années, les plages d'adresses IP de CCBot étaient documentées publiquement (par exemple, 38.107.191.66 - 38.107.191.119) (Source: <u>datadome.co</u>), mais maintenant le robot d'exploration est entièrement basé sur le cloud.

Robots.txt et éthique: Comme les robots d'exploration respectueux, CCBot respecte les règles robots.txt et les balises nofollow (Source: alchetron.com), il évite donc les pages explicitement interdites par les propriétaires de sites. Il se concentre sur le contenu accessible au public (pages HTML) et stocke le contenu brut des pages (HTML et texte) dans les archives d'exploration. Contrairement à l'Internet Archive, qui cherche à préserver les pages à des fins d'archivage et de relecture (y compris les images, les scripts et les comportements côté client) (Source: dallascard.github.io), l'accent de Common Crawl est mis sur le contenu textuel et les métadonnées utiles pour l'exploration de données et l'apprentissage automatique. Plus précisément, Common Crawl ne stocke ni n'analyse en détail les images, vidéos, CSS ou autres ressources statiques - l'accent est mis sur le texte HTML brut et les métadonnées associées. Cela rend le corpus de Common Crawl plus directement utile pour le TLN et l'analyse de données, au détriment d'un instantané visuel complet.

Méthodologie d'exploration : Common Crawl effectue généralement une **exploration d'un mois**, ce qui signifie qu'il exécute CCBot en continu pour récupérer des pages pendant environ un mois, puis publie les résultats sous forme d'« archive d'exploration ». Il répète cela environ tous les mois. Historiquement, le calendrier a varié : au cours des premières années, il y avait environ 4 explorations par an (Source: <u>alchetron.com</u>), mais plus tard, c'est devenu mensuel. Chaque exploration mensuelle part d'un vaste



ensemble d'URL de départ (points d'entrée initiaux) sur le web public et suit les liens pour découvrir de nouvelles URL, en élaguant en cours de route à l'aide d'heuristiques basées sur le domaine pour maintenir une large couverture. Le résultat de chaque exploration est une collection de fichiers WARC (archives compressées de pages récupérées) ainsi que les métadonnées associées (par exemple, tables d'URL, extraits de texte, graphes de liens) (Source: <u>alchetron.com</u>). Vers la mi-2012, Common Crawl a également commencé à publier le texte et les métadonnées extraits de chaque exploration, plutôt que de simples fichiers WARC bruts (Source: <u>alchetron.com</u>).

Échelle et croissance : L'ampleur des opérations de Common Crawl est massive. Selon une interview de 2023, chaque mois, Common Crawl collecte 3 à 5 milliards de pages web, soit « 500 fois plus de pages web que [l'ensemble de Wikipédia] » (Source: www.96layers.ai). Les données mensuelles compressées sont de l'ordre de 90 téraoctets (environ 400 To non compressés) (Source: www.96layers.ai). Sur plus d'une décennie, Common Crawl a accumulé des centaines de milliards de pages. Selon un rapport (avril 2024), il a été noté qu'« au cours de ses 17 ans d'histoire, Common Crawl a collecté plus de 250 milliards de pages web » (Source: www.96layers.ai). Sa propre page d'accueil (fin 2025) annonce « plus de 300 milliards de pages couvrant 15 ans » (Source: commoncrawl.org). (Ces chiffres sont globalement cohérents, compte tenu de la poursuite de l'exploration.) Pour situer, lors de son lancement début 2013, le jeu de données inaugural de Common Crawl comprenait environ 5 milliards de pages (≈81 téraoctets) (Source: nonprofitquarterly.org) (Source: www.thekurzweillibrary.com). Mi-2015, les explorations archivées couvraient environ 1,8 milliard de pages (145 To) sur 4 explorations annuelles (Source: alchetron.com). Aujourd'hui, l'exploration mensuelle dépasse à elle seule ces totaux antérieurs.

En plus du contenu des pages, Common Crawl publie également des **graphes de liens au niveau de l'hôte et du domaine** et d'autres jeux de données dérivés (par exemple, les URL contenant une requête donnée, ou des approximations de PageRank au niveau du domaine). Ceux-ci sont disponibles sur sa page *Données* et sur GitHub, et sont mis à jour régulièrement. Les archives WARC brutes et le texte traité sont hébergés dans **Amazon S3** (Jeu de données public AWS) et sur des sites miroirs. Les utilisateurs peuvent télécharger des explorations spécifiques par mois/année via HTTP ou utiliser des outils de mégadonnées (par exemple, Amazon Athena, Spark) pour interroger les données sur place. Common Crawl fournit également des outils d'aide et des index (par exemple, un index d'URL) pour faciliter la recherche de pages d'intérêt.

Dans l'ensemble, la technologie d'exploration de Common Crawl a évolué mais est restée ouverte. Elle utilise des composants standard et bien connus (Apache Nutch, cloud Amazon) et du code open source pour le traitement des données. En tant que projet à but non lucratif, il tire parti du cloud de manière créative : il évite de payer les coûts de stockage en restant sur le niveau gratuit d'AWS, et il contourne les frais de transmission de données (sortie) en encourageant l'analyse sur la plateforme AWS. L'infrastructure de base de Common Crawl est relativement simple, mais le résultat est énorme : des téraoctets de données web ouvertes agrégées et maintenues comme une ressource commune (Source: www.96layers.ai) (Source: dallascard.github.io).

Jeu de données et statistiques

Le jeu de données public de Common Crawl est l'un des plus grands corpus de texte existants, comparable en taille au stockage des principaux moteurs de recherche. Les statistiques clés concernant le corpus (à la mi-2025) sont :

- Taille du corpus : Plus de 300 milliards de pages web uniques (documents HTML) collectées (Source: commoncrawl.org). (À titre de comparaison, c'est des milliers de fois plus grand que l'ensemble de Wikipédia en anglais.)
- **Période couverte :** Instantanés mensuels de 2008 ou 2009 à aujourd'hui (plus de 15 ans) (Source: <u>commoncrawl.org</u>). Chaque instantané contient généralement les pages explorées ce mois-là. La collection s'enrichit chaque année.
- Taux de croissance mensuel : Généralement 3 à 5 milliards de pages par mois, produisant environ 90 To compressés
 (~400 To non compressés) chaque mois (Source: www.96layers.ai) (Source: commoncrawl.org). Sur une année, cela représente de l'ordre de 30 à 60 milliards de pages et des centaines de téraoctets.
- Fréquence d'exploration: Généralement une exploration par mois (bien qu'au début, c'était moins fréquent). L'archive est cumulative dans le sens où chaque exploration est un nouvel instantané, mais en pratique, les utilisateurs peuvent combiner des données sur plusieurs mois.
- Volume de données: Des centaines de téraoctets par exploration répartis dans des fichiers WARC, plus le texte dérivé et les métadonnées dans des fichiers adjacents. Par exemple, l'exploration inaugurale de 2013 était de 81 To (Source: nonprofitquarterly.org), et les explorations modernes sont plus importantes. Au total, les archives de Common Crawl représentent plusieurs pétaoctets de données compressées (le rapport 2024 de Mozilla cite « plus de 9,5 pétaoctets » de données Common Crawl) (Source: www.mozillafoundation.org).



Utilisation dans la littérature de recherche: Plus de 10 000 articles de recherche ont cité Common Crawl comme source de données (Source: commoncrawl.org) (Source: dallascard.github.io). Ce chiffre semble avoir à peu près doublé tous les quelques ans. (Le nombre exact est difficile à vérifier, mais le site web revendique fièrement « cité dans plus de 10 000 articles de recherche » (Source: commoncrawl.org), et des données indépendantes montrent que le nombre était bien inférieur en 2013.)

Ces chiffres approximatifs démontrent l'échelle massive des données. Il est à noter que seules quelques organisations privées (Google, Microsoft, Amazon, Facebook) disposent d'une capacité d'exploration web à une échelle comparable – et elles gardent les données propriétaires. En revanche, l'archive de Common Crawl est publiquement répertoriée sur <u>AWS Open Data</u> et d'autres miroirs, permettant à **quiconque** de la télécharger ou de l'analyser (Source: <u>registry.opendata.aws</u>).

Il est important de noter que Common Crawl précise que son jeu de données n'est **pas** le « web entier » et n'est pas garanti d'être complet. La couverture est biaisée en faveur des pages web accessibles en anglais (les sites bloqués via robots.txt sont exclus, et les grandes plateformes comme Facebook bloquent l'exploration). Une étude de Mozilla de 2024 a expressément averti que « *Traiter Common Crawl sans esprit critique comme une "copie du web" revient à déclarer qu'une sous-section relativement petite de pages web principalement en anglais est représentative du monde entier.* » (Source: www.mozillafoundation.org). En pratique, Common Crawl représente le « web visible » (la partie accessible à partir de liens HTML typiques) à la date de chaque exploration, en mettant l'accent sur la diversité (il ne se concentre pas exclusivement sur les domaines de premier niveau) et la fraîcheur.

Malgré ses limites, l'ampleur des données de Common Crawl les rend extrêmement précieuses. Il **dépasse de loin** tout jeu de données statique que la plupart des chercheurs pourraient collecter par eux-mêmes. Les modèles de langage naturel modernes utilisent couramment des **centaines de milliards de mots** provenant de Common Crawl. Par exemple, l'intégration de mots GloVe de Stanford (2014) a été entraînée sur **840 milliards de jetons** extraits de Common Crawl (Source: huggingface.co). Et les grands LLM ingèrent régulièrement des milliers de pages web informelles de Common Crawl (comme détaillé ci-dessous). Les données sont également utilisées dans l'analyse de graphes web, la recherche en récupération d'informations (par exemple, la construction de moteurs de recherche pour le jeu de données ClueWeb (Source: commoncrawl.org), et l'exploration spécifique à un domaine (comme l'extraction de texte parallèle pour la traduction automatique (Source: huggingface.co).

Le tableau 1 ci-dessous résume certaines de ces métriques et faits clés :



MÉTRIQUE/FAIT	VALEUR/DESCRIPTION	SOURCE
Année de fondation	2007 (établie en tant qu'organisation à but non lucratif 501(c)(3) en 2007) [9†L0 L24]	
Fondateur et Président	Gil Elbaz (technologue, co-fondateur d'Applied Semantics/AdSense) [47†L0-l	
Conseil consultatif (notable)	Peter Norvig de Google, Joi Ito du MIT, Nova Spivack, Carl Malamud	[30†L36-L38], [47†L19-L24], [45†L10-L18]
Type d'organisation	Organisation à but non lucratif 501(c)(3) (Californie)	[9†L0-L4], [7†L19- L24]
Âge/Période du jeu de données	2008/2009 - présent (plus de 15 ans de pages web mensuelles)	[9†L10-L17], [2†L20-L24]
Total de pages collectées	~300+ milliards de pages web (cumulatif) [9†L10 [2†L20	
Croissance mensuelle (pages)	~3 à 5 milliards de nouvelles pages ajoutées par mois (moyenne) [2†L:	
Taille des données mensuelles	~90 téraoctets compressés (~400 To non compressés) par exploration mensuelle	[2†L20-L24]
Critères d'inclusion	Pages HTML publiques (respectant robots.txt) ; accent sur le texte brut (pas d'images/vidéos).	[52†L22-L31], [19†L28-L31]
Utilisations notables du projet	Entraînement IA/ML (GPT-3, PaLM, etc.), intégrations de mots (jetons GloVe 840B), corpus de recherche (C4, The Pile), moteurs de recherche	[60†L23-L30], [61†L32-L39], [52†L49-L57]
Citations de recherche (approx.)	>10 000 articles publiés citant Common Crawl	[9†L12-L17], [52†L34-L40]
Jeu de données hébergé par Amazon	Hébergé via AWS Open Data (gratuit pour les utilisateurs via S3/Athena/AWS)	[19†L33-L39], [25†L12-L19] (registre AWS)
Plus grande couverture LLM	~80-85% des jetons d'entraînement de GPT-3 proviennent de Common Crawl (Source: www.mozillafoundation.org); ~64% des LLM étudiés (2019-2023) utilisent CC (Source: www.mozillafoundation.org).	

(Tableau 1 : Faits et statistiques clés sur Common Crawl, avec sources citées.)

Historique et développement

Le développement de Common Crawl peut être examiné chronologiquement à travers plusieurs étapes clés :



- 2007 Lancement du projet : Gil Elbaz « m'a approché avec une vision ambitieuse il voulait créer une exploration web ouverte et à but non lucratif » (Source: www.novaspivack.com). En 2007, il a officiellement fondé la Common Crawl Foundation. Parmi les premiers collaborateurs figuraient Nova Spivack et Carl Malamud, qui sont devenus membres du conseil d'administration (Source: www.novaspivack.com). À ce stade, seules quelques personnes y travaillaient (Elbaz lui-même, Ahad Rana en tant qu'ingénieur principal, quelques bénévoles). Spivack raconte : « Gil et l'ingénieur principal, Ahad Rana, se sont alors mis au travail pour construire la chose. » (Source: www.novaspivack.com). L'objectif était de créer « le premier index de recherche du Web véritablement ouvert, à but non lucratif, de 5 milliards de pages » (Source: www.novaspivack.com). (En effet, les premières données d'exploration publiées vers 2013 contenaient environ 5 milliards de pages, 81 To, comme rapporté par le MIT Tech Review (Source: nonprofitquarterly.org) (Source: www.thekurzweillibrary.com).)
- 2008-2011 Premières explorations: Après sa création, Common Crawl a commencé des explorations mensuelles (environ trimestrielles) d'une partie du web. Au cours de ces années, les volumes de données étaient plus petits; les premiers articles de blog indiquent seulement quelques téraoctets par exploration. L'accent était mis sur la construction du pipeline (explorateur basé sur Nutch, archives WARC, processus Hadoop simples pour extraire le texte). Initialement, l'équipe a écrit du code personnalisé, mais en 2013, elle a annoncé passer à Apache Nutch et adopter le format de fichier WARC pour toutes les données d'exploration (Source: alchetron.com). L'utilisation d'Amazon S3 pour le stockage a probablement commencé à cette époque.
- 2012 Partenariat avec Amazon AWS: Un tournant majeur s'est produit en 2012 lorsque Amazon Web Services a accepté Common Crawl dans son programme Public Datasets (Source: alchetron.com). AWS a accepté d'héberger les archives d'exploration dans son cloud sans frais. Cela a été crucial cela a permis à Common Crawl de passer des gigaoctets aux pétaoctets sans supporter les dépenses de stockage. (En parallèle, AWS et Common Crawl ont ensuite collaboré sur des concours; par exemple, AWS a offert aux participants 50 \$ de crédits pour utiliser les données (Source: commoncrawl.org).) Également fin 2012, la société de moteurs de recherche Blekko a fait don de métadonnées de ses propres explorations (févrieroctobre 2012) à Common Crawl (Source: alchetron.com). Les journaux de Blekko ont aidé à améliorer la couverture de l'exploration et à réduire les pages indésirables (spam, pornographie, manipulations SEO) (Source: alchetron.com).
- 2013 Lancement officiel et reconnaissance: Début 2013, la première grande publication publique de Common Crawl (l'« index de 5 milliards de pages ») a attiré l'attention des médias. Le MIT Technology Review (via le blog de Ray Kurzweil) a publié un article en janvier 2013 intitulé « Une base de données gratuite de l'ensemble du Web pourrait engendrer le prochain Google » (Source: www.thekurzweillibrary.com). L'article soulignait que « Common Crawl offre plus de cinq milliards de pages Web, disponibles gratuitement afin que les chercheurs et les entrepreneurs puissent tenter des choses autrement possibles uniquement pour ceux qui ont accès aux ressources de Google. » (Source: www.thekurzweillibrary.com). À cette époque, Peter Norvig et Joi Ito avaient rejoint le conseil consultatif (Source: www.thekurzweillibrary.com). Le propre site et tableau de bord de Common Crawl a été lancé, faisant la promotion de l'archive de données de dix ans et attirant les premiers utilisateurs de recherche.
- 2014-2019 Expansion des données et croissance de l'écosystème : Au milieu des années 2010, Common Crawl a poursuivi ses explorations mensuelles, et le jeu de données cumulatif a augmenté rapidement. Chaque année, davantage de recherche et développement a été construit sur ces données. Les événements importants incluent :
 - 2014-2015: Extraction de données structurées: Common Crawl a commencé à extraire du texte et des métadonnées des pages brutes et à les publier aux côtés des fichiers WARC. Des données pour des langues comme l'espagnol, l'allemand, etc. ont été mises à disposition. La communauté a également développé des outils pour interroger les données sur place, tels que Recipes et Index (via AWS Athena).
 - 2016: Introduction de CCBot v2.0 avec un agent utilisateur mis à jour (Source: datadome.co) et des améliorations concernant le respect de robots.txt. Le rôle de Common Crawl dans la recherche s'est consolidé, car des tâches de PNL comme GloVe (84 Go) ont utilisé les données CC (Source: huggingface.co).
- 2017-2019: Le jeu de données a dépassé les dizaines de milliards de pages. Pendant cette période, l'Europe a lancé le Norvig Web Data Science Award (soutenu par Common Crawl et SURFSara), encourageant l'utilisation académique des données. De plus, l'équipe d'ingénierie principale est restée petite ; lors d'entretiens, ils ont mentionné n'avoir que 3 employés environ en 2017 (Source: www.96layers.ai). En 2019, Common Crawl était reconnu comme une source clé pour l'entraînement des modèles neuronaux, bien qu'il restait encore peu connu du grand public.



- 2020-2022 Boom de l'IA: Le boom de l'IA de l'ère COVID a propulsé Common Crawl sous les feux de la rampe. GPT-3 d'OpenAl (publié mi-2020) a utilisé Common Crawl comme source de données principale. Des équipes de recherche derrière des modèles comme Grover (Zellers et al., 2019) se sont explicitement entraînées sur CC pour la génération de fausses nouvelles (Source: dallascard.github.io). RoBERTa (2019) de Meta et T5 de Google ont également puisé dans des corpus dérivés de CC. En 2020, les données de Common Crawl ont été intégrées à de grands jeux de données de recherche comme « C4 » (utilisé pour T5) et « The Pile » (un corpus anglais de 800 Go) tous deux reconnaissant publiquement CC comme un composant majeur (Source: dallascard.github.io). Le public a commencé à entendre parler de « trillions de tokens » extraits du web pour l'IA, et Common Crawl a été identifié comme une source clé. Cependant, Common Crawl lui-même est resté petit ; il a été rapporté qu'au moment du lancement de GPT-3, l'organisation n'avait peut-être qu'un seul employé y travaillant (Source: www.96layers.ai).
- 2023-2025 Ère actuelle et reconnaissance publique: En 2023 et 2024, Common Crawl a connu un regain d'attention publique dû à deux facteurs: (a) l'essor de l'IA générative, pour laquelle les données ouvertes de CC sont essentielles; et (b) les controverses juridiques concernant les contenus protégés par le droit d'auteur dans les données d'entraînement. Début 2024, la Fondation Mozilla a publié un rapport approfondi (basé sur des entretiens avec le personnel de Common Crawl) intitulé « Training Data for the Price of a Sandwich: Common Crawl's Impact on Generative AI. » (Source: www.mozillafoundation.org). Ce rapport a révélé des statistiques actualisées (9,5 Po de données, 84 % des tokens de GPT-3 provenant de CC) et a fourni des informations actualisées sur l'organisation. À peu près au même moment, une affaire juridique notable (New York Times contre OpenAI/Microsoft) a propulsé Common Crawl à la une, car le contenu du NYT avait été extrait par CC et ainsi utilisé par inadvertance dans GPT-3 (Source: www.mozillafoundation.org). L'équipe de Common Crawl a également annoncé de nouveaux services (par exemple, l'hébergement d'un Common Crawl Index interrogeable (Source: commoncrawl.org) et un engagement communautaire accru (articles, tutoriels, hackathons).

Tout au long de son histoire, Common Crawl est resté fidèle à sa mission originale d'accès ouvert. Il n'est jamais devenu un moteur de recherche commercial ou un fournisseur de données. Au lieu de cela, il s'est concentré sur la construction d'un pipeline robuste et évolutif et d'une communauté autour des données ouvertes. La direction du projet souligne régulièrement que « fournir des données d'entraînement pour l'IA n'a jamais été le but principal de Common Crawl », et qu'ils ont toujours accueilli une large base d'utilisateurs (les chercheurs en IA n'étant qu'un groupe) (Source: www.96layers.ai). Néanmoins, comme nous le verrons, l'avènement de l'IA générative a rendu Common Crawl plus influent que jamais – à la fois pour le bien (permettre la recherche) et pour la controverse (préoccupations concernant le droit d'auteur et les biais).

Détails techniques des données de Common Crawl

Formats de données et accès

Chaque exploration de Common Crawl produit un ensemble de fichiers au format **WARC** (Web ARChive), qui regroupe des séquences de réponses HTTP (les pages web récupérées) avec des métadonnées. Ces fichiers WARC sont le résultat brut de l'exploration, généralement nommés par la date et l'identifiant de l'exploration. En plus des WARC, Common Crawl publie une variété de fichiers d'accompagnement :

- **Texte extrait (fichiers WAT) :** Pour chaque WARC, un fichier « WAT » correspondant contient des métadonnées analysées (par exemple, en-têtes HTTP, liens, métadonnées JSON).
- Texte extrait (fichiers WET): Un fichier « WET » diffuse le texte brut extrait de chaque page HTML (essentiellement le contenu textuel nettoyé). Ceux-ci permettent aux utilisateurs d'analyser rapidement le texte sans avoir à analyser le HTML euxmêmes
- Index d'URL (CDX): Un index CSV/JSON de toutes les URL récupérées et de leurs décalages dans les WARC, utile pour interroger des sites ou des pages spécifiques.
- Graphes web: Données graphiques reliant des pages ou des domaines (par exemple, des graphes de liens d'hôte à hôte).
 Ceux-ci sont fournis périodiquement (par exemple, annuellement) pour étudier la connectivité.
- Tables de domaines : Fichiers agrégés listant tous les domaines explorés et le nombre de pages.



Tous ces fichiers sont stockés dans des **buckets AWS S3** (et mis en miroir ailleurs). Common Crawl encourage l'utilisation de l'analyse dans le cloud (par exemple, Amazon Athena ou EMR) pour interroger les données à grande échelle. Par exemple, Amazon Athena permet des requêtes SQL sur l'index de toutes les URL ou même sur le contenu WARC s'il est correctement structuré. Le coût d'exécution de ces requêtes est faible (et parfois couvert par des crédits), ce qui rend pratique pour les équipes de recherche d'extraire des jeux de données de Common Crawl sans copier des téraoctets sur leurs serveurs locaux.

Common Crawl lui-même fournit des outils de développement et de la documentation (par exemple, le projet « Index to WARC Files and URLs » (Source: registry.opendata.aws). Mais il existe aussi un écosystème externe dynamique : de nombreux projets GitHub et tutoriels (par exemple, CC-pyspark, commoncrawljob) aident les nouveaux utilisateurs à démarrer. La liste de diffusion publique de Common Crawl et les communautés Slack/Discord sont actives avec des astuces et du code partagé.

CCBot (Robot d'exploration de Common Crawl)

Le **robot d'exploration web** lui-même, surnommé **CCBot**, fonctionne en continu lors de chaque exploration mensuelle. Il fonctionne à peu près comme suit : un ordonnanceur maître distribue des instances de robot d'exploration (sur AWS EC2) qui récupèrent des pages en parallèle, en suivant la liste des URL à visiter. De nouvelles URL sont ajoutées à la file d'attente à mesure que des liens sont découverts. Le robot d'exploration utilise les fonctionnalités standard de Nutch : respect de robots.txt, limitation automatique par domaine, et logique de déduplication pour éviter d'explorer sans fin le même contenu (par exemple, suppression des paramètres de session).

CCBot s'identifie avec une chaîne d'agent utilisateur, mais Common Crawl recommande aux webmasters de ne pas se fier uniquement à cela pour la mise en liste blanche, car des robots d'exploration malveillants peuvent l'usurper (Source: datadome.co). (Au lieu de cela, les propriétaires de sites peuvent utiliser les plages d'adresses IP AWS connues pour identifier le trafic de CCBot.) Bien qu'étant un utilisateur légitime, les adresses IP de CCBot proviennent de pools AWS dynamiques, de sorte que certains sites le bloquent ou le limitent par inadvertance. Common Crawl s'efforce d'être un robot d'exploration « poli ». Par exemple, il fait tourner les plages d'adresses IP, se retire des sites surchargés et autorise certaines erreurs d'exploration. Les administrateurs de serveurs qui souhaitent respecter les normes de la communauté peuvent explicitement autoriser CCBot en ajustant leur robots.txt (Common Crawl fournit de la documentation sur la façon de procéder).

Au fil du temps, CCBot a été affiné pour l'efficacité. L'architecture actuelle (en 2025) utilise un système distribué et tolérant aux pannes sur AWS, coordonné par l'équipe principale (dirigée par un « ingénieur d'exploration »). L'exploration de mai 2025, par exemple, a couvert **2,47 milliards de pages** (voir le rapport du sommet Twitter (Source: commoncrawl.org). Au total, le système s'est avéré **évolutif**: Common Crawl note fièrement que son exploration est désormais « gargantuesque », bien au-delà de la capacité de tout chercheur universitaire à la dupliquer (Source: nonprofitquarterly.org).

Pipeline de traitement des données

Les pages brutes explorées subissent un pipeline de traitement avant leur publication. Les étapes clés comprennent :

- Extraction de liens : Identifier tous les hyperliens sur chaque page à ajouter à la frontière d'exploration. Construire des graphes de liens (au niveau du domaine et de l'hôte) pour l'analyse.
- Déduplication de contenu: Filtrer les pages identiques ou quasi-identiques pour réduire le gaspillage et les biais. Common Crawl applique une déduplication agressive au niveau du document et de la page afin que les données archivées aient une redondance minimale.
- Extraction de texte : Supprimer le HTML/CSS et extraire le contenu textuel, qui est stocké dans les fichiers « WET ». Cela inclut la détection de la langue (Common Crawl se concentre généralement sur le texte anglais mais capturera également d'autres langues).
- Métadonnées HTTP: Enregistrer les en-têtes de réponse, le type de contenu et les informations du serveur pour chaque récupération (dans les fichiers WAT).
- Gestion des erreurs: Enregistrer toutes les erreurs de récupération ou les délais d'attente dans un fichier « errata ».
 Common Crawl maintient un journal d'errata qui liste les URL ou les domaines qui échouent constamment, afin d'améliorer les explorations futures.



Le résultat final est un produit de données riche : pour un mois donné, un utilisateur peut récupérer non seulement les blobs HTML bruts, mais aussi un corpus parallèle de phrases (le texte WET) et toute la structure des hyperliens. Le code du pipeline est open source, et des améliorations (par exemple, une meilleure analyse HTML, gestion de JavaScript) sont périodiquement intégrées.

(En février 2023, Common Crawl a annoncé sur son blog son intention d'expérimenter le *pré-rendu* des pages nécessitant JavaScript – mais fin 2025, le corpus principal reste centré sur le HTML.)

Caractéristiques du jeu de données

- Distribution linguistique: Les menus de Common Crawl révèlent que le jeu de données est multilingue, mais fortement orienté vers l'anglais. Selon le rapport de Mozilla, l'exploration est « principalement en anglais » avec une couverture régionale variable. Par exemple, des jeux de données de 50 millions d'articles de presse allemands (Source: commoncrawl.org) et d'autres corpus spécifiques à des langues ont été dérivés de CC, mais l'exploration brute contient beaucoup plus de contenu anglais.
- **Diversité des sites :** Common Crawl essaie d'équilibrer l'étendue et la profondeur. Il inclut les sites majeurs (actualités, ecommerce, blogs) ainsi que les sites à longue traîne. Cependant, il ne cible pas le « web profond » ni les pages protégées par mot de passe. Il ne peut pas non plus explorer les sites qui interdisent les robots ou nécessitent des identifiants.
- Instantanés temporels: Chaque exploration mensuelle est horodatée. Par conséquent, les archives de Common Crawl
 peuvent être utilisées pour étudier l'évolution du web (par exemple, comment une page ou un domaine change au fil du
 temps). Cependant, Common Crawl n'est pas une archive continue comme la Wayback Machine il ne préserve pas chaque
 version d'une page quotidiennement; il fournit principalement une « capture » par URL par mois (à moins que la page ne
 change et ne soit réexplorée plus tard).

Dans l'ensemble, les données de Common Crawl sont extrêmement volumineuses et assez représentatives du web public (sous réserve des robots et de l'accès). C'est *la* plus grande archive web publiquement disponible pour la recherche, combinant volume et accessibilité.

Cas d'utilisation et impact

Le jeu de données ouvert de Common Crawl a permis une vaste gamme d'applications. Nous organisons son utilisation en plusieurs grandes catégories :

1. IA et apprentissage automatique (LLM, Embeddings, etc.)

Common Crawl est devenu la source de données fondamentale pour le traitement du langage naturel et l'IA à grande échelle. Pratiquement tous les modèles linguistiques modernes ont puisé dans ces données. Par exemple :

- **GPT-3 et ChatGPT :** Lorsque OpenAl a entraîné GPT-3 (qui est à la base de ChatGPT), la majorité de ses tokens d'entraînement provenaient de Common Crawl. L'article publié par OpenAl sur GPT-3 montre que « *la plus grande quantité de données d'entraînement provient de Common Crawl »* (Source: <u>datadome.co</u>). Une analyse de Mozilla corrobore cela : elle a constaté que *plus de 80 % des tokens de GPT-3* provenaient de Common Crawl (Source: <u>www.mozillafoundation.org</u>). (Les GPU s'entraînent généralement sur plusieurs corpus ; pour GPT-3, les autres sources étaient WebText2, des livres et Wikipédia. Mais Common Crawl représentait la plus grande partie.) Parce que GPT-3 alimente directement les chatbots et les assistants IA, le contenu de Common Crawl (bon ou mauvais) « parle » essentiellement aux utilisateurs finaux via l'IA.
- Autres grands modèles linguistiques: De nombreux autres LLM notables ont été construits sur les données de CC:
 - Les modèles T5 et basés sur BERT de Google ont incorporé des sous-ensembles de Common Crawl.
 - RoBERTa de Facebook a été entraîné sur un mélange de données CC et d'actualités en 2019.
 - Des modèles open source comme GPT-NeoX d'EleutherAI et des modèles plus petits tels que GPT-2 ont utilisé CC.
 - Le modèle **Grover** (2019) de Zellers *et al.* un modèle pour générer et détecter les fausses nouvelles a explicitement utilisé Common Crawl pour le texte web (Source: <u>dallascard.github.io</u>).



- Plus récemment, la plupart des nouveaux modèles (Bellatrix, LLaMA, etc.) utilisent des pipelines comme The Pile ou RefinedWeb, qui à leur tour sont tirés des instantanés de Common Crawl (Source: <u>dallascard.github.io</u>). En effet, les instantanés de Common Crawl sont reconditionnés dans des jeux de données dérivés (par exemple, C4, Colossal Clean Crawls) qui alimentent les charges de travail d'entraînement à grande échelle.
- Une enquête sur 47 LLM divers (2019-2023) a révélé qu'« au moins 64 % » d'entre eux avaient été entraînés sur des données de Common Crawl (Source: www.mozillafoundation.org). Cela inclut les modèles de nouvelle génération comme ChatGPT-4 (via GPT-4), LLaMA de Meta, Mistral, Claude 2, etc. (Certains modèles peuvent également utiliser des données propriétaires ou mixtes, mais CC reste un pilier.)
- Embeddings de mots et outils PNL: Le jeu de données a permis des ressources PNL fondamentales. Les embeddings GloVe classiques (840 milliards de tokens, anglais) et les embeddings FastText (600 milliards de tokens) sont tous deux entraînés sur le texte de CC (Source: huggingface.co). Des corpus open source comme Colossal Clean Crawls (C4) et des jeux de données multilingues dérivés de Common Crawl alimentent les modèles de traduction et les résumeurs. La recherche en modélisation de sujets, analyse de sentiments, récupération d'informations, et plus encore utilise souvent CC comme source de texte brut. Par exemple, une étude de 2019 a construit un corpus parallèle bilingue à partir de CC pour la traduction automatique (Source: huggingface.co).
- Chatbots et assistants IA: Au-delà de l'entraînement de modèles hors ligne, certains services effectuent une exploration en temps réel de CC pour prendre en charge l'IA. Par exemple, DeepSeek et certaines plateformes de recherche « basées sur l'IA » ingèrent des pages CC pour fournir leurs réponses. De nombreux bots IA s'appuient également sur CC pour vérifier les faits ou augmenter les réponses, car c'est un index pratique du web public.
- Données pour les modèles de vision et multimodaux : Bien que Common Crawl contienne principalement du texte, il contient également des URL d'images (et parfois des métadonnées d'images). Des entreprises comme TinEye exploitent l'index d'URL d'images de CC pour créer des services de recherche d'images inversée (Source: nonprofitquarterly.org). (TinEye a explicitement utilisé Common Crawl pour trouver des images similaires à une image de requête.) Certains modèles de vision IA utilisent des légendes textuelles alignées sur CC ou du texte alternatif dans les données de CC pour les associer à des images.

En résumé, **les chercheurs et les entreprises en IA utilisent massivement Common Crawl** comme source de données gratuite. Son omniprésence dans l'entraînement des modèles a soulevé à la fois des opportunités (faire progresser l'IA) et des préoccupations (biais, droit d'auteur) – plus de détails ci-dessous.

2. Recherche académique et scientifique

Le corpus Common Crawl est largement cité dans la recherche académique, toutes disciplines confondues :

- Langage naturel et science du web: Les chercheurs analysent l'utilisation et les modèles linguistiques. Par exemple, CC a
 été utilisé pour étudier la structure des hyperliens (qui lie qui sur le web), géolocaliser des actualités (un jeu de données de 50
 millions d'articles de presse allemands a été construit à partir de CC (Source: commoncrawl.org), et analyser la lisibilité ou les
 expressions courantes sur le web. Les travaux sur les graphes web (théorie des graphes appliquée aux domaines) utilisent
 souvent les données de graphes de liens de CC (Source: commoncrawl.org).
- Exploration de données et analyse de Big Data: Le jeu de données illustre les « grandes données ouvertes ». Les chercheurs testent des algorithmes d'exploration de texte à grande échelle (clustering, détection d'anomalies, analyse thématique) sur CC. La capacité d'accéder à des pétaoctets de données réelles a permis des études comparatives de pipelines de traitement de texte.
- Études en récupération d'information (RI) : Common Crawl est utilisé pour construire des moteurs de recherche expérimentaux. Par exemple, Elastic ChatNoir à Bauhaus Weimar est conçu pour rechercher dans les archives ClueWeb et Common Crawl (Source: commoncrawl.org). Les chercheurs en RI évaluent également les algorithmes de classement sur des sous-ensembles de CC, ou utilisent CC comme référence pour le contenu des pages web. L'équipe de Common Crawl ellemême fournit une API « Simple Speedy Search » (CCSS) pour des recherches rapides par mots-clés sur l'index.
- Cybersécurité et mesure des abus: La nature à grande échelle de CC permet de rechercher des modèles malveillants. Par
 exemple, l'article « Lurking Malice in the Cloud » (ACM 2016) a scanné toutes les pages de CC pour trouver des scripts intégrés
 liés à des domaines de logiciels malveillants connus (Source: huggingface.co). Les chercheurs ont utilisé CC pour quantifier la



prévalence des en-têtes HTTP (in)sécurisés, des bibliothèques obsolètes ou des scripts de cryptojacking sur les sites web populaires.

- Économie et sciences sociales: Les spécialistes des sciences sociales utilisent CC comme un indicateur du discours public. Par exemple, une étude a utilisé CC pour analyser la modération de contenu et la censure; la recherche du Citizen Lab « Banned Books » a analysé des pages de produits Amazon extraites via CC pour détecter les politiques de censure (Source: commoncrawl.org). D'autres cas d'utilisation incluent le suivi de la désinformation en matière de santé, l'analyse de la propagande politique ou l'étude de la diffusion de contenu en plusieurs langues sur le web ouvert.
- Indices de citation et cartographie de la science: La disponibilité de milliards de citations savantes glanées dans les textes de CC a même permis la méta-recherche. Par exemple, la réexécution d'analyses de citations et la construction de graphes de connaissances à une échelle colossale.

Il est à noter que le site web de Common Crawl lui-même met en avant de nombreux articles de recherche : il regroupe des liens vers des travaux publiés exploitant les données de CC (Source: commoncrawl.org). Les citations couvrent NeurIPS/ICLR pour le PNL, les conférences WWW/WWW pour l'analyse web, et des revues dans les domaines de l'IA, des sciences de l'information et des sciences sociales computationnelles.

3. Applications commerciales et industrielles

Au-delà du monde universitaire, de nombreuses entreprises et startups ont bâti des produits basés sur les données de Common Crawl. Voici quelques exemples notables :

- Recherche d'images TinEye: Comme mentionné, TinEye (par Idée Inc.) utilise Common Crawl pour indexer les images.
 Lorsqu'un utilisateur soumet une image, TinEye la hache et effectue une recherche sur les données d'images collectées à partir de CC pour trouver des images similaires (Source: nonprofitquarterly.org). CC a fourni une source d'images et de leurs URL vaste et gratuite, permettant à TinEye de lancer une entreprise viable sans avoir à explorer le web par elle-même.
- Analyse d'impact Lucky Oyster: Lucky Oyster Labs (acquise par Rendever) a utilisé Common Crawl pour l'écoute sociale
 et l'analyse des tendances. Ils ont développé des outils sur CC pour « comprendre ce dont les gens discutent sur le web » en
 tant que moteur d'analyse (Source: nonprofitquarterly.org). (L'article du NPQ mentionne Lucky Oyster comme une startup
 exploitant CC, bien que les détails soient désormais rares.)
- Search-as-a-Service Cas Crate.IO: Certaines entreprises ont développé des connecteurs et des moteurs pour interroger les données de CC. Par exemple, Crate.IO a publié un blog sur « l'importation à partir de sources de données personnalisées » à l'aide d'un plugin, montrant comment alimenter les archives de CC dans leur base de données SQL (Source: commoncrawl.org). De même, « CommonCrawlJob » et « CommonCrawlScalaTools » sont des projets GitHub qui aident à charger les données de CC dans des systèmes de big data. Il s'agit principalement de preuves de concept ou d'outils de développement.
- Moteurs de recherche de startups: Au moins une équipe entrepreneuriale (Elastic ChatNoir (Source: commoncrawl.org) a
 construit une interface de moteur de recherche spécifiquement pour les clones Common Crawl du jeu de données ClueWeb.
 Une autre, les instantanés web ouverts de Carrot Search, a expérimenté avec CC. Il y a un intérêt à créer des moteurs de
 recherche à but non lucratif ou alternatifs utilisant CC comme backend de données évitant ainsi la nécessité d'explorer le web
 par soi-même.
- Marketing et SEO: Certaines entreprises d'analyse SEO utilisent CC pour estimer l'accès au site ou l'analyse des concurrents.
 Bien que la plupart des produits SEO commerciaux s'appuient sur des crawlers propriétaires, CC offre un pool de données gratuit pour évaluer le nombre de pages globales ou les tendances de contenu. Par exemple, les lignes de code des outils SEO comme Majestic ou Ahrefs pourraient incorporer des données CC pour l'analyse des backlinks, bien que les détails soient généralement propriétaires.
- Publicité et Business Intelligence: Les entreprises de données (y compris Factual, la société fondée par Gil Elbaz) ont intégré les données de CC pour enrichir les ensembles de données commerciales. Par exemple, le nombre de domaines, la fraîcheur des sites et la classification du contenu peuvent être glanés à partir de CC pour alimenter le ciblage publicitaire ou les outils de marketing B2B. Cependant, en raison de la nature automatisée des données, les informations basées sur CC doivent être validées avec soin pour un usage commercial.



Le tableau 2 (ci-dessous) résume quelques cas d'utilisation et projets illustratifs qui exploitent les données de Common Crawl :

UTILISATEUR/PROJET	CAS D'UTILISATION	SOURCE / NOTES
TinEye	Recherche d'images inversée (trouver des images similaires en explorant)	Utilise des images explorées par CC (Source: nonprofitquarterly.org). (IDée Inc.)
Lucky Oyster	Analyse des tendances sociales/culturelles	Startup utilisant CC pour analyser les tendances de contenu web (Source: <u>nonprofitquarterly.org</u>).
GloVe (Stanford)	Embeddings de vecteurs de mots (840 milliards de tokens de CC)	CC a fourni le texte pour le modèle GloVe (Source: huggingface.co).
GPT-3/ChatGPT	Données d'entraînement pour grand modèle linguistique (~80% des tokens de CC)	Rapport Mozilla : « <i>Plus de 80 % des tokens de GPT-3</i> provenaient de Common Crawl. » (Source: www.mozillafoundation.org).
Modèles linguistiques	Entraînement/affinage (RoBERTa, T5, LLaMA, etc.)	Les LLM (2019-2023) utilisent souvent des corpus basés sur CC (Source: dallascard.github.io) (Source: www.mozillafoundation.org).
Moteurs de recherche	Construction d'index de recherche alternatifs (ex. ChatNoir)	Elastic ChatNoir : recherche de données CC (Source: commoncrawl.org). (Bauhaus-Weimar)
Recherche PNL	Analyse statistique de texte web (modèles thématiques, résumé)	Des dizaines d'articles universitaires dans les domaines du PNL citent CC.
Métriques web	Études sur la censure/liberté d'expression (ex. censure Amazon)	Citizen Lab « Banned Books » a utilisé CC (Source: commoncrawl.org) ; autres articles de science web.

(Tableau 2 : Exemples sélectionnés de l'utilisation des données de Common Crawl en pratique, avec citations.)

En plus de ces exemples, le propre site web de Common Crawl répertorie de **nombreux projets** : des ensembles de données ouverts (WikiSQL à partir de tables web), des expériences de recherche basées sur le cloud, des tutoriels Elasticsearch et des cours universitaires, tous construits sur les données de CC (Source: <u>commoncrawl.org</u>). De manière anecdotique, Gil Elbaz a commenté que « **si vous n'êtes pas Google, OpenAl ou Microsoft, presque tout le monde s'appuie sur Common Crawl** » pour les données à grande échelle (Source: <u>www.96layers.ai</u>). Cela souligne à quel point CC est devenu omniprésent pour toute organisation qui ne peut pas déployer son propre crawler web à l'échelle de Google.

Études de cas

Pour illustrer plus concrètement l'impact de Common Crawl, nous décrivons deux études de cas détaillées : l'une sur l'IA/l'entraînement de modèles et l'autre sur la recherche ouverte.

Étude de cas : GPT-3 et la révolution des LLM

Prenons l'exemple très médiatisé de GPT-3 (2020) d'OpenAI et de ses modèles apparentés. Ces « Transformers génératifs préentraînés » atteignent des capacités impressionnantes en langage naturel, mais leur puissance dérive de vastes données d'entraînement. Common Crawl a joué un rôle de premier plan :

Composition de l'ensemble de données: L'article sur GPT-3 (Brown et al. 2020) énumère les sources de données:
 WebText2 (le propre crawl d'OpenAl de pages liées à Reddit), Google Books, Wikipedia et Common Crawl. En termes de taille brute, Common Crawl était de loin le plus grand. Une analyse ultérieure confirme que « la plus grande quantité de données



d'entraînement provient de Common Crawl » (Source: <u>datadome.co</u>). Le rapport de Mozilla précise que **plus de 80** % **de tous les tokens utilisés par GPT-3 provenaient de CC** (Source: <u>www.mozillafoundation.org</u>).

- Modèle résultant: GPT-3-175B, avec 175 milliards de paramètres, a été entraîné sur 570 Go de données textuelles filtrées (environ 500 milliards de tokens). Si 80 % provenaient de CC, cela signifie environ 456 Go de texte CC. Cette échelle serait impossible sans un corpus web existant. La disponibilité de CC a permis à OpenAl de ne pas avoir à allouer de ressources pour explorer le web par elle-même à ce moment-là (bien qu'ils aient probablement aussi eu des données web internes).
- Utilisation professionnelle: Lors du lancement de GPT-3, il a été rapidement intégré dans des produits (par exemple, Copilot de Microsoft, ChatGPT d'OpenAl en 2022). Ces services agissent alors comme une « couche d'IA » au-dessus de CC. Certains utilisateurs craignent que, lorsque ChatGPT fournit des réponses, il ne régurgite du texte provenant de pages Common Crawl sans attribution. En effet, le rapport de Mozilla note que les modèles basés sur CC produisent souvent du contenu biaisé ou protégé par le droit d'auteur car ils ont tendance à mémoriser les données d'entraînement.
- Implications légales (Affaire NYT): Fin 2023, The New York Times a poursuivi OpenAI, alléguant que les données d'entraînement de ChatGPT (GPT-3.5/GPT-4) incluaient indûment du contenu du Times. Common Crawl est devenu une pièce maîtresse de la preuve car les articles du Times avaient été explorés dans CC avant l'entraînement du modèle, et OpenAI a utilisé ces instantanés de CC. Une fiche d'information de Mozilla explique: « Le contenu du NYT représentait une proportion significative des données de Common Crawl au moment où OpenAI a lancé ChatGPT, et constituait donc probablement une partie significative des données d'entraînement de GPT-3 » (Source: www.mozillafoundation.org). Cela souligne comment l'ouverture de CC peut involontairement entraîner une exposition juridique lorsque du texte protégé par le droit d'auteur est redistribué dans des modèles.
- **Diversité et biais :** Parce que tant de LLM s'appuient sur CC, les directives apprises dans CC se propagent largement. Si CC manque de contenu suffisant provenant de certaines langues ou données démographiques, les modèles peuvent sousperformer sur ces sujets. La recherche de Mozilla avertit que « *l'ensemble de données de Common Crawl inclut délibérément du contenu problématique (toxicité, discours de haine, etc.) afin de soutenir la recherche sur ces phénomènes.* » En revanche, de nombreux pipelines d'entraînement d'IA filtrent fortement CC (par exemple, ne conservent que les « pages anglaises de haute qualité ») (Source: www.mozillafoundation.org), ce qui signifie que la toxicité brute de CC peut influencer le comportement du modèle si elle n'est pas soigneusement supprimée.

En résumé, le cas GPT-3 montre que **Common Crawl est devenu l'épine dorsale de la recherche en lA générative** dans les années 2020. Il a considérablement réduit la barrière à l'entraînement de grands modèles. Le fait que les données d'une petite organisation à but non lucratif alimentent des systèmes d'IA de plusieurs millions de dollars est remarquable. Cela force également une remise en question : lorsqu'un ensemble de données ouvert alimente une IA à code source fermé, qui est responsable du contenu ? La direction de Common Crawl souligne que les données étaient destinées à toutes sortes d'analyses (y compris la recherche sur les discours de haine), et non explicitement à l'entraînement de modèles valant des milliards de dollars (Source: www.96layers.ai). Le débat communautaire tourne désormais autour de la manière de garantir que les modèles basés sur CC sont « fiables » (suppression des biais, respect du droit d'auteur, etc.) (Source: www.mozillafoundation.org).

Étude de cas : Recherche ouverte via Common Crawl

Un autre cas illustratif est la tentative de **construire des moteurs de recherche en utilisant les données de Common Crawl**. Alors que les entreprises expertes en web comme Google ou Bing développent leurs propres crawlers, certains groupes indépendants ont exploré l'utilisation de CC comme source de données pour des services de recherche alternatifs.

- Elastic ChatNoir: Des chercheurs de l'Université Bauhaus ont créé ChatNoir, une interface de recherche ouverte pour les corpus ClueWeb et CC (Source: commoncrawl.org). Ceci est destiné à la recherche en humanités numériques et en récupération d'informations. ChatNoir indexe les pages de Common Crawl et fournit une interface de recherche simple, permettant aux utilisateurs d'interroger l'archive de CC comme s'il s'agissait d'un moteur de recherche. Cela démontre qu'en principe, on peut utiliser CC comme « backend » pour la recherche.
- CC Search (Bêta): Common Crawl lui-même a lancé CC Search (maintenant opéré par l'équipe Creative Commons/WordPress) qui permet aux utilisateurs de rechercher par mots-clés dans CC. Le site web de CC mentionne des mises à jour comme « Gros changements pour CC Search Beta » fin 2024 (rédigé par Paola Villarrela). L'objectif est de rendre les données de CC plus accessibles (par exemple, en ajoutant la recherche par licence, langue, etc.).



- Propositions de startups: L'idée d'un « moteur de recherche à but non lucratif » a été évoquée périodiquement (même sur Hacker News (Source: dallascard.github.io). Même le titre de l'article du Nonprofit Quarterly était « Rencontrez Common Crawl, l'organisation à but non lucratif qui pourrait remodeler le web » (Source: nonprofitquarterly.org). Pour l'instant, Common Crawl lui-même reste uniquement axé sur les données (pas de portail de recherche utilisateur), mais des tiers peuvent s'en servir. L'existence de CC signifie que tout groupe bien doté en ressources pourrait lancer un moteur de recherche sans explorer le web par lui-même.
- Considérations pratiques : Il est important de noter que les données de Common Crawl ont des limitations pour la recherche : elles n'incluent pas le PageRank, les données de clics des utilisateurs ou une fraîcheur à jour au-delà de la granularité mensuelle. Certains sites web excluent CC, et l'ensemble de données est « gelé » à des points mensuels. Ainsi, un moteur basé sur CC serait partiellement obsolète. Néanmoins, des projets de recherche à petite échelle « spécifiques à un domaine » ont utilisé CC avec succès. Par exemple, une équipe de recherche pourrait restreindre CC aux domaines d'actualités et construire une recherche d'actualités spécialisée.

Dans le commerce électronique ou le SEO, certaines entreprises explorent CC pour collecter des informations ouvertes sur les données de produits ou les classements de sites. Il est rapporté qu'un blogueur (Claus Matzinger de Crate.IO) a écrit sur l'importation de données CC dans une base de données optimisée pour la recherche (Source: commoncrawl.org). Comme l'a dit un observateur de longue date de CC : « Si vous n'êtes pas Google, OpenAI ou Microsoft... presque tout le monde s'appuie sur Common Crawl » (Source: www.96layers.ai) pour au moins certaines données à grande échelle.

Ces cas montrent que Common Crawl a permis de **nouveaux types de services** que seuls les géants de la recherche pouvaient auparavant envisager. Bien qu'aucun moteur de recherche commercial majeur (avec des requêtes en direct) n'ait entièrement adopté CC, le projet a effectivement abaissé la barrière : construire un système de recherche expérimental ou académique sur Common Crawl est simple et rentable.

Analyse des données et résultats de recherche

Au-delà des anecdotes d'utilisation, les chercheurs ont analysé quantitativement Common Crawl lui-même. Voici quelques résultats représentatifs :

- Échelle des données: Une interview de 2024 avec Stefan Baack, chercheur chez Mozilla, a résumé le volume mensuel et historique de Common Crawl (Source: www.96layers.ai). Par exemple, il note que chaque archive mensuelle pèse 90 To compressés et que Common Crawl a accumulé « plus de 250 milliards de pages web » en 17 ans (Source: www.96layers.ai). Ces chiffres sont cohérents avec l'affirmation officielle du site de « plus de 300 milliards de pages » (Source: commoncrawl.org). Une telle analyse souligne la taille inégalée de CC.
- Métriques de citation: En explorant Google Scholar ou les bases de données bibliographiques, le personnel de Common Crawl a constaté que leurs données avaient plus de 10 000 citations dans la littérature universitaire (Source: commoncrawl.org) (Source: dallascard.github.io). Cela démontre la large adoption dans divers domaines. Les chercheurs ont indiqué que CC est utilisé dans des domaines aussi variés que la détection de spam web, les bibliothèques numériques, le journalisme (suivi des fausses nouvelles) et même l'informatique de la santé (par exemple, la recherche de désinformation médicale).
- Couverture linguistique et de sites: Le rapport de Mozilla souligne que l'anglais domine Common Crawl. Il montre le nombre de pages web par pays/langue et note que de nombreuses pages chinoises, japonaises et de médias sociaux (par exemple, Facebook, Twitter) sont manquantes ou sous-représentées en raison des restrictions d'exploration (Source: www.mozillafoundation.org). En fait, les pages des sites qui bloquent explicitement les crawlers sont absentes. Le rapport souligne également que l'objectif de CC de soutenir la « recherche sur les discours de haine » signifie qu'il inclut intentionnellement un tel contenu (Source: www.mozillafoundation.org), ce qui est un choix de conception (laissé non filtré pour permettre l'analyse). Cependant, ceux qui s'intéressent à l'entraînement des LLM filtrent souvent ces pages.
- Robustesse technique: L'analyse des données de journalisation de CC a été effectuée pour évaluer le processus de crawling web lui-même. Par exemple, l'article de Springer « Web Crawl Refusals: Insights from Common Crawl » a étudié comment les serveurs web bloquent ou ralentissent les crawlers, en utilisant les propres journaux de récupération de CC (Source: commoncrawl.org). Les résultats ont permis d'établir les meilleures pratiques pour le crawling (par exemple, comment gérer les faux blocages de type « fake chatgpt-bot »).



Richesse sémantique des données: Certains projets ont tenté d'annoter CC à grande échelle. Par exemple, en créant des graphes de connaissances en extrayant des entités et des relations du texte de CC. Le projet <u>CSRankings</u> de Stanford utilise CC pour évaluer la taille des publications CS de CVPR, ICML, NeurIPS (bien que ce soit une digression). Mais plus pertinent: des chercheurs ont utilisé CC pour construire des graphes de connaissances ouverts de « bon sens » en analysant des milliards de phrases.

En résumé, la **méta-analyse** de Common Crawl confirme son ampleur et son influence. Des études indépendantes ont validé les statistiques brutes du site et exploré ses biais. Ces études contribuent à l'amélioration du jeu de données (par exemple, en mettant en évidence les régions du web sous-crawlé) et à guider les utilisateurs sur l'utilisation appropriée (par exemple, en avertissant des problèmes de droits d'auteur) (Source: www.mozillafoundation.org).

Défis, limitations et problèmes

Bien que les données de Common Crawl soient puissantes, elles ne sont pas sans défis ni critiques :

- Biais et représentativité : Comme indiqué, CC est biaisé en termes de langue (principalement l'anglais) et de région (davantage États-Unis/UE). Certains domaines (comme le contenu africain et asiatique) sont sous-représentés. Cela peut biaiser toute analyse ou IA entraînée sur CC. Le rapport de Mozilla avertit explicitement que CC ne doit pas être traité comme un « substitut de l'ensemble du web » (Source: www.mozillafoundation.org). Les chercheurs complètent souvent CC avec d'autres corpus pour une meilleure couverture (par exemple, actualités, archives gouvernementales, collections spécifiques à une langue).
- Qualité du contenu : Common Crawl inclut délibérément une grande variété de contenus, ce qui signifie qu'il capture également des pages web de mauvaise qualité, spammées ou toxiques. Il n'y a pas de filtrage strict du « bon » par rapport au « mauvais » contenu par défaut. Pour certains cas d'utilisation (recherche linguistique, détection de biais), cette inclusivité est une caractéristique. Mais pour l'entraînement de l'IA, cela nécessite un nettoyage supplémentaire. Par exemple, l'article pertinent d'Ablestacks sur le Pile et des jeux de données similaires inclut plusieurs filtres pour supprimer les blasphèmes, le contenu pour adultes, le texte non-anglais, etc. L'analyse de Mozilla souligne que les développeurs d'IA doivent « éliminer » le contenu indésirable de CC si leur objectif est un entraînement de modèle sûr (Source: www.mozillafoundation.org). En pratique, de nombreuses pipelines d'IA (Aleph, Redwood, etc.) utilisent des listes participatives ou heuristiques pour filtrer CC.
- Droits d'auteur et licences: Les « Conditions d'utilisation » de CC stipulent que les pages web sont collectées sans égard aux droits d'auteur, en supposant que le texte sur le web public peut être utilisé (similaire au fonctionnement de Googlebot). Cependant, l'essor de l'IA a soulevé des questions juridiques. Le procès susmentionné du New York Times suggère que CC pourrait avoir aspiré des milliers d'articles protégés par le droit d'auteur sur des sites d'actualités, lesquels se sont ensuite retrouvés dans les paramètres de GPT-3. Cela illustre une tension: Common Crawl estime que sa collecte de données est légalement protégée (par exemple, en vertu des exceptions du DMCA pour la mise en cache/le crawling, ou sous l'idée d'« utilisation transformatrice » dans l'entraînement de l'IA). Mais les titulaires de droits ne sont pas d'accord. Common Crawl n'a pas spécifiquement demandé la permission à chaque créateur de contenu sur le web; il s'appuie fondamentalement sur les conditions d'utilisation d'Internet et le fichier robots.txt. Fin 2023, Common Crawl a clarifié qu'une fois le contenu sur CC, il est « là pour que tous l'utilisent (ce qui inclut le fine-tuning et l'inférence / l'augmentation par récupération) » (Source: www.mozillafoundation.org). Cette position est controversée.
- Comité et gouvernance : Comme CC est géré par des bénévoles, son avenir dépend de la bonne volonté continue et du soutien des sponsors. Il n'y a pas de financement garanti ni de dotation importante. Si les principaux donateurs technologiques retiraient leur soutien, les opérations de CC pourraient être compromises. Cependant, à partir de 2025, l'intérêt pour la conservation des projets de données web ouvertes semble élevé, compte tenu de l'intérêt législatif pour la régulation de l'IA et la science ouverte. Common Crawl a des plans (selon les dernières déclarations) pour diversifier son financement et éventuellement ajouter de nouvelles fonctionnalités (comme des métadonnées de licence, des API d'opt-out, etc.) pour répondre aux préoccupations des propriétaires de contenu.
- Limitations techniques: Le jeu de données est massif, mais il peut encore manquer du contenu généré dynamiquement ou caché derrière des formulaires. Les sites utilisant un rendu côté client intensif ou nécessitant JavaScript peuvent être partiellement invisibles pour les crawlers textuels uniquement. Certaines pages modernes (par exemple, les applications à page unique) avec peu de HTML statique pourraient ne pas être bien capturées. Common Crawl a expérimenté les navigateurs sans



tête (headless browsers), mais cela est coûteux. Par conséquent, CC peut sous-indexer les sites très modernes et fortement basés sur JavaScript. De plus, comme il effectue un passage par mois, il peut manquer des mises à jour rapides ou des pages éphémères. Les utilisateurs ayant besoin de données fraîches en temps réel ne peuvent pas se fier uniquement à CC.

Dans l'ensemble, l'équipe de Common Crawl reconnaît ces défis. Leur stratégie a été la transparence : ils publient fréquemment des articles de blog et des réponses pour expliquer la portée et les limites du jeu de données (par exemple, « Web Archiving File Formats Explained » (Source: commoncrawl.org). Ils encouragent les utilisateurs à considérer CC comme une infrastructure partagée, semblable à une expérience ouverte, plutôt qu'un produit parfait.

Orientations futures et implications

Pour l'avenir, Common Crawl se situe à l'intersection de plusieurs tendances en science des données et en gouvernance d'Internet :

- Amélioration de la qualité des données: Common Crawl pourrait adopter un filtrage ou un étiquetage plus avancé pour
 mieux servir les utilisateurs. Par exemple, la génération d'un sous-ensemble « nettoyé » du crawl (supprimant le spam probable
 ou le contenu pour adultes) pourrait favoriser une adoption plus large. Inversement, la création de sous-crawls spécialisés (par
 exemple, un crawl multilinque ou un crawl anglais de haute qualité) pourrait attirer de nouveaux publics.
- Propriétaires de contenu et permissions: À mesure que les débats sur les droits des données évoluent, Common Crawl pourrait mettre en œuvre des mécanismes d'opt-out. Déjà, certains sites proposent des règles DDD/Robots.txt pour l'exclusion de l'IA. Common Crawl s'est porté volontaire pour respecter les x-robot-tags bloquant tout crawling non-bot (style DRM). Les futurs systèmes pourraient permettre aux propriétaires de sites de demander la suppression de l'archive de CC. D'autre part, de tels opt-outs menacent l'uniformité des jeux de données pour les chercheurs. Le projet continuera probablement à collaborer avec des experts juridiques pour trouver un équilibre.
- Initiatives de recherche ouverte: Il y a une défense croissante de l'« infrastructure de recherche en tant que service public
 ». Common Crawl pourrait devenir la base de données d'une nouvelle génération de moteurs de recherche ouverts ou de
 graphes de connaissances. Par exemple, des projets comme OpenWebIndex (un projet proposé financé par l'UE) font écho à
 la mission de Common Crawl. Nous pourrions voir des partenariats où le crawl de Common Crawl alimente des index spécialisés
 (par exemple, un moteur de recherche académique de contenu éducatif, ou une recherche ouverte de shopping). La publication
 de l'API d'Index de Common Crawl (annoncée en 2023) montre un mouvement dans cette direction.
- IA et utilisation responsable : Étant donné que les données de Common Crawl alimentent l'IA générative, la fondation pourrait investir dans des fonctionnalités d'« éthique de l'IA ». Cela pourrait inclure des annotations (marquant les pages de propagande ou de désinformation sanitaire) ou l'intégration de filtres de débiaisage. Le rapport de Mozilla suggère que les développeurs devraient ajouter des « filtres de données robustes » (Source: www.mozillafoundation.org) ; Common Crawl pourrait lui-même commencer à proposer des versions pré-filtrées ou des outils de filtrage (par exemple, un filtre de toxicité).
- Analyses supplémentaires par Common Crawl: La fondation pourrait produire davantage d'analyses de données en interne. Par exemple, leur GitHub présente des tableaux de bord « Crawl Stats » et « Graph Stats ». L'extension de ceux-ci pour afficher des répartitions linguistiques en temps réel, des métriques de diversité de domaine ou des tendances sémantiques pourrait être précieuse. Cela aiderait à la fois les utilisateurs et les bailleurs de fonds à comprendre la portée de la ressource.
- Partenariats mondiaux: Pour améliorer la couverture, Common Crawl pourrait s'associer à des universités internationales ou des ONG pour alimenter le crawl avec davantage de contenu mondial (par exemple, via les 100 premiers domaines spécifiques à chaque pays). Il pourrait également collaborer avec des bibliothèques nationales (comme Europeana ou des archives web nationales) pour intégrer les « jardins clos » du web.

Plus largement, l'impact de Common Crawl suggère que les **biens communs de données** (infrastructure de données ouvertes) pourraient être un modèle viable pour d'autres domaines : imaginez des corpus ouverts d'articles scientifiques, d'images ou de capteurs environnementaux. Le succès de Common Crawl fournit un modèle : équipe minimale, sponsors cloud, données ouvertes. Il montre que, dans les bonnes conditions, « les données sont la nouvelle infrastructure publique ».

Conclusion



Common Crawl est né de la vision de Gil Elbaz d'un index web ouvert, et en près de deux décennies, il est devenu une ressource essentielle pour l'innovation basée sur les données. Son **histoire** est celle de débuts modestes (une petite organisation à but non lucratif en 2007) qui, grâce aux efforts de la communauté et au soutien du cloud, est devenue une **archive web gigantesque** (Source: nonprofitquarterly.org) (Source: www.96layers.ai). Il est né d'un engagement envers les données ouvertes et a adhéré à ce principe : rendre l'information à l'échelle du web démocratiquement accessible, et non propriétaire.

Aujourd'hui, Common Crawl est utilisé par des milliers de chercheurs et de développeurs dans le monde entier. Il alimente la pointe de l'IA (pratiquement tous les grands modèles linguistiques s'y fient) et permet à des start-ups qui, autrement, ne pourraient pas se permettre l'infrastructure de Google. Le tableau 2 de ce rapport a illustré quelques exemples concrets, mais un décompte exhaustif serait encore plus long. Sa présence dans plus de **10 000** publications académiques (Source: commoncrawl.org) (Source: dallascard.github.io) témoigne de son influence.

Cependant, une grande puissance s'accompagne de responsabilités et de complications. L'utilisation de Common Crawl dans l'entraînement de l'IA a soulevé des questions sociales et juridiques – d'autant plus que les modèles génératifs façonnent le discours public. L'équipe de Common Crawl en est consciente et s'est engagée avec la communauté sur la manière d'utiliser les données de manière responsable. Le rapport de Mozilla et d'autres analyses suggèrent que CC fera partie des débats sur l'éthique de l'IA et les droits d'auteur pour les années à venir (Source: www.mozillafoundation.org). (Source: www.mozillafoundation.org).

Pour l'avenir, la trajectoire de Common Crawl semble orientée vers une expansion continue et une intégration plus profonde avec la recherche ouverte. À mesure que la puissance de calcul augmente et que l'IA recherche toujours plus de données, la valeur de l'archive web ouverte de Common Crawl augmentera probablement. La communauté autour de lui pourrait s'étendre, passant peut-être d'une petite équipe à un consortium collaboratif plus vaste. Des projets naissants visent à étendre ses capacités (tels que des index de recherche plus riches ou des options de filtrage) qui pourraient façonner l'ère du « Search 2.0 » (Source: commoncrawl.org).

En somme, l'histoire complète de Common Crawl est une étude de cas sur la façon dont une petite initiative bien ciblée peut considérablement ouvrir les biens communs de données. Il a commencé comme une réponse aux craintes de monopole dans la recherche web, et il a effectivement ouvert des portes à l'innovation. Son fondateur Gil Elbaz et ses collaborateurs ont réussi à créer « le web comme une base de données géante », accessible à tous (Source: nonprofitquarterly.org). L'histoire de Common Crawl – du premier crawl de cinq milliards de pages à des milliers de milliards de pages aujourd'hui – illustre la puissance de l'infrastructure ouverte. Son rôle futur s'approfondira probablement à mesure que la société fera face aux avantages et aux défis de l'IA à l'échelle du web et de la science ouverte.

Toutes les affirmations ci-dessus sont étayées par des sources citées provenant de la propre documentation de Common Crawl, de rapports médiatiques, d'entretiens et d'analyses savantes (Source: commoncrawl.org) (Source: www.96layers.ai) (Source: <a href="www.96layers.ai) (Sou

Étiquettes: exploration-commune, exploration-web, donnees-entrainement-llm, donnees-ouvertes, gil-elbaz, big-data, depot-web, apache-nutch

AVERTISSEMENT

Ce document est fourni à titre informatif uniquement. Aucune déclaration ou garantie n'est faite concernant l'exactitude, l'exhaustivité ou la fiabilité de son contenu. Toute utilisation de ces informations est à vos propres risques. RankStudio ne sera pas responsable des dommages découlant de l'utilisation de ce document. Ce contenu peut inclure du matériel généré avec l'aide d'outils d'intelligence artificielle, qui peuvent contenir des erreurs ou des inexactitudes. Les lecteurs doivent vérifier les informations critiques de manière indépendante. Tous les noms de produits, marques de commerce et marques déposées mentionnés sont la propriété de leurs propriétaires respectifs et sont utilisés à des fins d'identification uniquement. L'utilisation de ces noms n'implique pas l'approbation. Ce document ne constitue pas un conseil professionnel ou juridique. Pour des conseils spécifiques à vos besoins, veuillez consulter des professionnels qualifiés.