Rankstudie

Pourquoi Cloudflare bloque les robots d'IA par défaut : Une analyse

By RankStudio Publié le 9 octobre 2025 29 min de lecture



Résumé Analytique

L'émergence de l'IA générative a bouleversé le modèle symbiotique traditionnel entre les éditeurs de contenu et les robots d'exploration web. Historiquement, les moteurs de recherche comme Google **exploraient les sites web** pour améliorer l'expérience de recherche, redirigeant le trafic utilisateur vers la source originale. En revanche, les systèmes d'IA modernes (par exemple, <u>ChatGPT</u>, Gemini, Claude) déploient des **robots d'IA** avancés qui collectent le contenu web pour entraîner de grands modèles linguistiques, souvent sans rediriger les utilisateurs vers la source. Ce changement a suscité une vive inquiétude chez les éditeurs, qui voient leurs revenus publicitaires et d'abonnements diminuer tandis que les entreprises d'IA profitent du contenu librement collecté.

Cloudflare, un fournisseur de CDN et d'infrastructure internet de premier plan (protégeant environ 20 % d'Internet (Source: www.windowscentral.com), a réagi à ce changement de paradigme en apportant des modifications importantes à sa politique. Mi-2025, Cloudflare a inversé sa position sur l'exploration par l'IA: plutôt que d'autoriser (optionnellement) les robots d'exploration par défaut, il bloque les robots d'IA par défaut sur les nouveaux sites web. Les propriétaires de sites web peuvent toujours choisir d'autoriser des robots d'exploration spécifiques, mais seulement après avoir donné une permission explicite et clarifié l'intention du robot (entraînement, inférence ou recherche) (Source: www.infosecurity-magazine.com) (Source: adquilly.me). Cette mesure a été accompagnée d'une suite de nouveaux outils - robots.txt géré, signaux de contenu et un système de "paiement par exploration" - conçus pour donner aux éditeurs le contrôle sur leurs données.

Le raisonnement principal de Cloudflare est de protéger les intérêts économiques des créateurs de contenu et de préserver un web libre et ouvert à l'ère de l'IA. La direction de Cloudflare soutient que sans changement, l'incitation à produire du contenu original disparaîtra. Comme l'ont averti Page et le co-fondateur *Matthew Prince*, l'exploration incontrôlée par l'IA "prive

les créateurs de contenu de revenus" et menace l'avenir d'Internet (Source: <u>adgully.me</u>). En imposant un modèle basé sur la permission et un opt-out par défaut pour le scraping par l'IA, Cloudflare vise à **rétablir l'équilibre** sur le web : les propriétaires de sites retrouvent leur autonomie (et une compensation potentielle) sur leur contenu (Source: <u>adgully.me</u>) (Source: <u>adgully.me</u>).

Ce rapport fournit une analyse complète de la nouvelle politique de blocage par défaut de Cloudflare, examinant le contexte technique (robots.txt et l'exploration), l'économie du contenu en évolution, les données et outils de Cloudflare, les réactions de l'industrie, des études de cas et les implications futures. Nous rassemblons des données sur l'activité des robots d'exploration, citons des opinions d'experts et des déclarations de l'industrie, et considérons de multiples perspectives (éditeurs, développeurs d'IA, régulateurs) pour expliquer **pourquoi Cloudflare a agi de la sorte**, et ce que cela augure pour le web.

Introduction et Contexte

L'architecture ouverte d'Internet a historiquement permis aux moteurs de recherche d'explorer et d'indexer le contenu, au bénéfice des utilisateurs et des propriétaires de sites. Le **Robots.txt**, introduit en 1994 (Source: blog.cloudflare.com) (Source: www.arrayanalytics.io), a permis aux webmasters de donner des instructions de base aux robots d'exploration sur ce qu'il fallait indexer ou éviter. Les bots conformes (notamment Googlebot) obéissaient à ces directives, générant du trafic vers les sites via les résultats de recherche. Pendant des décennies, cela a créé une **situation gagnant-gagnant**: les éditeurs gagnaient en visibilité et en revenus publicitaires, tandis que les entreprises de recherche construisaient de meilleurs services.

Cependant, l'essor des grands modèles linguistiques a perturbé cet équilibre. Les entreprises d'IA (par exemple, OpenAl, Google, Anthropic, Meta) déploient des **robots d'exploration web sophistiqués** (souvent appelés *bots IA*, *spiders IA* ou *scrapers IA*) pour collecter des ensembles de données massifs directement depuis le web. Contrairement aux robots d'exploration de recherche traditionnels, ces agents IA ne renvoient pas nécessairement les utilisateurs à la source. Au lieu de cela, ils utilisent le contenu scrapé pour générer des réponses dans des applications propriétaires ou pour entraîner des modèles. Les utilisateurs s'appuient de plus en plus sur des <u>résumés ou réponses générés par l'IA</u> (par exemple, ChatGPT ou les aperçus IA de Google) au lieu de cliquer sur les sites web originaux.

Cela a de profondes implications pour les créateurs de contenu en ligne. Sans trafic entrant, les vues publicitaires et l'intérêt des abonnés peuvent diminuer, sapant l'incitation économique à produire du contenu de qualité. Les éditeurs ont observé des **baisses spectaculaires** du trafic de référence provenant des moteurs de recherche, attribuées aux systèmes d'IA qui fournissent des "réponses" sans renvoyer de liens. Comme l'a noté le PDG de Cloudflare lors d'un sommet à Cannes, il y a dix ans, Google explorait environ 2 pages pour chaque visiteur envoyé à un éditeur ; aujourd'hui, les utilisateurs "suivent moins de notes de bas de page", réduisant drastiquement l'engagement avec le matériel source (Source: www.axios.com). Avec les robots d'IA, le déséquilibre est bien plus aigu : les données de Cloudflare montrent des **ratios d'exploration IA par rapport aux visites** de l'ordre de milliers, dépassant de loin les niveaux modestes des moteurs de recherche (Source: blog.cloudflare.com) (Source: blog.cloudflare.com) (voir Tableau 1).

Tableau 1 : Ratios d'exploration par rapport aux références pour les robots d'exploration web (juin 2025) (Source: blog.cloudflare.com). En termes simples, un ratio d'exploration par rapport aux références de X:1 signifie X visites par un robot d'exploration pour un clic de référence vers le site.

BOT/PLATEFORME	RATIO D'EXPLORATION PAR RAPPORT AUX RÉFÉRENCES
Google Search	~14:1
OpenAl (ChatGPT/GPTBot)	~1 700 : 1
Anthropic (ClaudeBot)	~73 000 : 1

Comme l'illustre le Tableau 1, les robots d'exploration d'entraînement IA visitent les sites *des ordres de grandeur plus* par référence que Google. En termes pratiques, une entreprise d'IA comme OpenAl pourrait demander **1 700 pages** à un site pour chaque visite d'utilisateur que ce site reçoit via les réponses de ChatGPT (Source: <u>blog.cloudflare.com</u>) (Source: <u>blog.cloudflare.com</u>). Pour Anthropic, l'écart est encore plus grand (rapporté à ~73 000:1). En revanche, le modèle classique de Google était d'environ une douzaine d'explorations par visite (Source: <u>blog.cloudflare.com</u>) (Source: <u>blog.cloudflare.com</u>).

Cette asymétrie extrême des données rompt le modèle "exploration contre trafic". Les éditeurs craignent désormais que les clients de l'IA puissent consommer leur contenu à grande échelle sans crédit ni compensation. Dans certains cas, les systèmes d'IA présentent même le contenu directement dans les résultats de recherche (par exemple, les extraits d'IA de Google), érodant davantage les clics vers les articles originaux. Les analyses des sociétés de licence de contenu et les poursuites judiciaires (par exemple, les poursuites du New York Times, de Ziff Davis contre OpenAl (Source: apnews.com) (Source: www.reuters.com) soulignent la perception des éditeurs d'une menace existentielle. Dans ce contexte, de nombreux éditeurs et défenseurs ont appelé à des contrôles plus stricts, y compris le respect de robots.txt ou le blocage pur et simple du scraping non autorisé (Source: www.reuters.com) (Source: www.reuters.com). Cloudflare, compte tenu de sa position privilégiée en tant que proxy et fournisseur de gestion de bots pour des millions de sites, a suivi de près ces tendances. En réponse, ils ont introduit de nouvelles fonctionnalités et des politiques par défaut pour aider les propriétaires de sites à reprendre le contrôle de leur contenu. Les sections à venir analyseront ce que Cloudflare a fait et pourquoi – situant leurs actions dans le contexte historique et technique plus large de l'exploration web et des droits de contenu.

Contexte Historique : Robots.txt et l'Exploration Web

Le **Protocole d'Exclusion des Robots**, incarné par le fichier robots.txt à la racine d'un site web, a été formalisé au milieu des années 1990 (initialement comme une convention informelle) pour aider les propriétaires de sites à guider les bots de recherche. Un robots.txt peut inclure des directives telles que Disallow ou Allow, spécifiant quels user-agents (bots) peuvent accéder à quelles parties du site (Source: blog.cloudflare.com) (Source: www.arrayanalytics.io). Il est crucial de noter que la conformité à robots.txt est volontaire: les robots d'exploration sont censés le respecter par courtoisie, et non en vertu d'une règle exécutoire (Source: blog.cloudflare.com) (Source: www.arrayanalytics.io). Les premiers bots majeurs (Googlebot, Bingbot, etc.) ont dûment respecté ces règles, permettant une interaction transparente : les sites web pouvaient bloquer les explorations indésirables sans cacher le contenu aux utilisateurs humains.

Au fil du temps, l'utilisation de robots.txt est devenue une pratique courante parmi les sites. Les données de Cloudflare montrent qu'environ un tiers des principaux domaines avaient un robots.txt à la mi-2025 (Source: blog.cloudflare.com). Cependant, même lorsqu'il était présent, peu de sites l'avaient explicitement configuré pour bloquer les robots d'exploration liés à l'IA. Les données Radar de Cloudflare indiquaient qu'à la mi-2025, seulement environ 7,8 % des principaux sites interdisaient nommément le "GPTBot" d'OpenAI, et des fractions encore plus petites bloquaient des bots comme anthropic-ai ou ClaudeBot (Source: blog.cloudflare.com). En d'autres termes, la plupart des créateurs de contenu n'avaient pas pleinement utilisé robots.txt pour exprimer leurs préférences concernant l'IA.

Pendant ce temps, de nombreux robots d'exploration modernes **ignorent ou contournent** robots.txt. Le problème est devenu urgent : Reuters a rapporté que « diverses entreprises d'IA contournent le Protocole d'Exclusion des Robots (robots.txt) pour extraire du contenu des sites d'éditeurs » (Source: www.reuters.com). Par exemple, le moteur de recherche IA **Perplexity** a été accusé par Cloudflare/d'autres de scraping malgré des règles Disallow explicites (Source: www.itpro.com) (Source: www.reuters.com) (Source: <a

En résumé, robots.txt a commencé comme une humble courtoisie standard du web, mais sa nature volontaire limite son application à l'ère de l'IA. Ce contexte explique la motivation de Cloudflare à aller plus loin : associer les signaux robots.txt à des blocages plus forts, appliqués au niveau du réseau, et à des politiques par défaut qui ne dépendent pas des propriétaires de sites les embauchant explicitement.

L'Ascension des Robots d'IA et la Rupture de l'Échange de Contenu

Historiquement, les spécialistes du SEO et les créateurs de contenu considéraient les robots d'exploration comme des alliés. Les spiders de Google rendaient le contenu de grande valeur découvrable, augmentant les vues de pages et les revenus publicitaires. Cette symbiose se fragmente désormais. Les applications d'IA modernes servent souvent des réponses directes ou des résumés aux utilisateurs, leur donnant ce dont ils ont besoin sans nécessiter un clic de retour vers le site web original (Source: adgully.me). La logique financière du web est ainsi compromise : un rapport de Reuters de 2025 a noté la baisse spectaculaire du **trafic de clics pour l'accès** alors que les résumés générés par l'IA supplantent les liens de recherche (Source: www.reuters.com) (Source: www.reuters.com).

Les analyses de trafic internes de Cloudflare rendent cela évident. À la mi-2025, l'équipe Radar de Cloudflare a rapporté que Google fournissait environ 14 requêtes d'exploration par visite de référence, tandis que les propres robots d'exploration d'OpenAl demandaient environ 1 700 pages par référence, et ceux d'Anthropic environ 73 000 (Source: blog.cloudflare.com) (Source: blog.cloudflare.com). Ce déséquilibre massif signifie que le contenu est extrait à grande échelle sans trafic correspondant. Cloudflare explique que cela "rompt clairement la relation 'exploration en échange de trafic' qui existait auparavant entre les robots d'exploration de recherche et les éditeurs" (Source: blog.cloudflare.com).

L'aspect axé sur les données de la décision de Cloudflare est clair : les éditeurs ne reçoivent plus les avantages de l'ouverture. Comme l'a formulé une analyse, les robots d'IA sont des "robots gourmands en données [qui extraient] du contenu créé par l'homme sans permission et sans le payer" (Source: www.infosecurity-magazine.com). En l'absence de visiteurs entrants, les sites ne génèrent aucune impression publicitaire et manquent des abonnements potentiels. Les grandes entreprises de contenu (par exemple, Condé Nast, Gannett, USA Today Network) ont publiquement soutenu les mesures de Cloudflare, citant explicitement les pertes de revenus et l'utilisation gratuite inéquitable du contenu comme motivation (Source: adgully.me) (Source: www.reuters.com). Cloudflare lui-même a fait écho à ce sentiment : il a averti que sans rééquilibrage, "l'avenir d'Internet est en péril" car les créateurs perdent leur incitation (Source: adgully.me).

En somme, l'appétit de l'IA pour les données a mis sous pression les modèles de revenus traditionnels. L'adoption par Cloudflare du blocage par défaut des bots est une réaction directe à ces pressions économiques. En contrôlant l'accès des robots d'exploration au niveau de la couche réseau, Cloudflare et ses clients visent à réintroduire le quid pro quo du web ouvert.

Données de Cloudflare et Résultats des Projets Pilotes

Au-delà des reportages externes, Cloudflare a accumulé ses propres preuves du problème de l'exploration par l'IA. Dans un billet de blog de 2025, l'entreprise a présenté des statistiques détaillées sur le trafic de bots vers les sites protégés par Cloudflare (Source: blog.cloudflare.com) (Source: adgully.me). Les principales conclusions sont les suivantes :

- Dominance des nouveaux bots d'IA: À la mi-2025, le GPTBot d'OpenAl était devenu le bot le plus répandu sur les sites Cloudflare, dépassant les crawlers traditionnels comme Googlebot et d'autres bots de grandes entreprises technologiques (Source: blog.cloudflare.com). Par exemple, les requêtes de GPTBot avaient même dépassé celles du crawler d'Amazon (voir le graphique dans [10]).
- Chute de la part d'exploration non-GPTAI: La part des sites accédés par les anciens scrapers (comme Bytespider de ByteDance) a chuté après les premières tentatives de blocage de Cloudflare. À partir de juillet 2024, la part d'accès de Bytespider a diminué d'environ 71 %, tandis que de nombreuses requêtes étaient explicitement bloquées par les paramètres du site (Source: blog.cloudflare.com).
- Adoption généralisée du blocage: Plus d'un million de sites sur Cloudflare ont activé activement la fonction de blocage des "scrapers d'IA" en un clic, introduite en juillet 2024 (Source: <u>blog.cloudflare.com</u>) (Source: <u>adgully.me</u>). Cela démontre un fort désir de blocage de la part des éditeurs. (En fait, Cloudflare a noté que cette adoption était l'impulsion pour faire du blocage la valeur par défaut pour les nouveaux sites (Source: <u>www.infosecurity-magazine.com</u>) (Source: <u>adgully.me</u>).)
- Sous-utilisation de robots.txt: Seulement environ 37 % des principaux domaines disposaient d'un fichier robots.txt. Parmi ceux-ci, très peu répertoriaient les crawlers d'IA dans les règles Disallow. Par exemple, en juillet 2025, seulement environ 7,8 % des principaux sites interdisaient GPTBot, et moins de 5 % interdisaient d'autres bots d'IA majeurs (Source: blog.cloudflare.com). Ces lacunes ont montré à Cloudflare que la gestion manuelle de robots.txt ne suivait pas le rythme des nouvelles menaces de bots.

Ces données renforcent la raison pour laquelle Cloudflare est intervenu. Les chercheurs de Cloudflare ont explicitement conclu que la plupart des sites web ne limitaient pas de manière proactive l'accès à l'IA, soit par ignorance, soit par manque de ressources techniques. En proposant des solutions gérées, Cloudflare pouvait combler cette lacune.

Parallèlement, les données réseau de Cloudflare montrent une **activité explosive des crawlers d'IA**. Dans un rapport, l'équipe Radar de Cloudflare a constaté une forte augmentation de l'exploration globale par les bots de recherche/assistants IA (par exemple, une hausse de 18 % d'un mois à l'autre au début de 2025 (Source: <u>noise.getoto.net</u>). Même si les volumes de requêtes individuelles peuvent être faibles par bot, l'agrégat est énorme en raison de la flotte de bots des startups d'IA qui se développe rapidement (Source: <u>workmind.ai</u>) (Source: <u>workmind.ai</u>). Cloudflare note que l'*infrastructure* requise pour servir ces crawlers - serveurs, bande passante - impose des coûts aux hébergeurs web, de sorte que le scraping non réglementé nuit également aux performances du site (Source: <u>workmind.ai</u>).

Collectivement, ces analyses ont conduit Cloudflare à croire qu'il disposait à la fois d'un **argument de vente technique** et d'une **justification éthique** pour le blocage par défaut des bots. Les données ont apporté un soutien quantitatif aux plaintes anecdotiques des éditeurs et ont éclairé le réglage fin des nouvelles fonctionnalités.

Les nouveaux outils de contrôle du contenu IA de Cloudflare

Pour résoudre le problème de l'exploration, Cloudflare a déployé plusieurs outils, aboutissant à la nouvelle politique de blocage par défaut. Ces initiatives peuvent être résumées comme suit :

FONCTIONNALITÉ/POLITIQUE	DESCRIPTION	DATE DE LANCEMENT
Blocage IA en un clic	Un interrupteur configurable par l'utilisateur (gratuit sur tous les plans) pour bloquer <i>toutes</i> les chaînes d'agent utilisateur de crawlers IA connus. Cela arrête immédiatement de nombreux bots IA à la périphérie du réseau.	Juillet 2024 (Source: adgully.me)
robots.txt géré avec signaux de contenu	Un service automatisé où Cloudflare crée ou met à jour le robots.txt du site pour inclure des directives spécifiques à l'IA (par exemple, interdire l'entraînement d'IA). Étend également le fichier avec de nouvelles balises d'utilisation de l'IA (ai-train, ai-input, etc.) afin que les propriétaires puissent déclarer comment le contenu de leur site peut être utilisé (Source: www.cloudflare.net).	Juillet 2025 (Source: www.cloudflare.net)
Blocage IA par défaut à l'inscription	Les nouveaux domaines ajoutés à Cloudflare sont désormais interrogés s'ils souhaitent autoriser les crawlers d'IA. La réponse par défaut est non , installant des règles robots.txt qui interdisent ou bloquent les bots d'IA. Les propriétaires de sites peuvent ultérieurement choisir d'autoriser des crawlers spécifiques (Source: adgully.me) (Source: adgully.me). De cette façon, chaque nouveau site démarre dans un état "sûr".	Juillet 2025 (Source: adgully.me)
Audit des crawlers IA et blocage granulaire	Outils de tableau de bord et d'API pour identifier précisément quels crawlers visitent un site, et les bloquer ou les autoriser sélectivement. Cloudflare a introduit des analyses granulaires du trafic de bots et des modèles en un clic pour bloquer des agents utilisateurs de bots IA spécifiques (Source: blog.cloudflare.com) (Source: adgully.me).	Septembre 2024 (Source: <u>adgully.me</u>)
Paiement par exploration (Bêta)	Un mécanisme permettant aux propriétaires de contenu de facturer les entreprises d'IA pour l'exploration. Les opérateurs de sites peuvent exiger un paiement (signalé par HTTP 402) pour les bots qui souhaitent accéder à du contenu au-delà des autorisations standard (Source: www.reuters.com). En effet, cela permet des négociations ou des licences concernant l'utilisation des données.	Juillet 2025 (bêta) (Source: www.reuters.com)

Tableau 2 : Résumé des initiatives de contrôle du contenu IA de Cloudflare (2024-2025). Les dates correspondent au lancement des versions bêta ou aux annonces des fonctionnalités.

Ces fonctionnalités reflètent un passage à un modèle basé sur la permission. Auparavant, les crawlers bénéficiaient d'un consentement implicite selon l'éthique du "web public" (sauf blocage manuel). Désormais, Cloudflare instaure un paradigme d'opt-in : les bots doivent être explicitement autorisés. Par exemple, comme l'a formulé *Stephanie Cohen* (CSO de Cloudflare),

dans le nouveau système, "les entreprises d'IA devront désormais obtenir une permission explicite pour accéder au contenu, y compris en précisant si leur intention est l'entraînement, l'inférence ou la recherche" (Source: www.infosecurity-magazine.com).

Le lancement d'un blocage par défaut sur les nouveaux sites est un élément clé de ce changement. En interrogeant les propriétaires de sites dès le départ et en bloquant par défaut, Cloudflare rend la politique applicable. Une explication officielle a noté que demander à chaque nouveau client lors de la configuration "élimine la nécessité pour les propriétaires de pages web de configurer manuellement leurs paramètres pour se désinscrire" (Source: adgully.me). En pratique, cela signifie qu'immédiatement après l'activation de Cloudflare, le contenu d'un nouveau domaine est (par défaut) protégé des bots IA. Le propriétaire doit prendre des mesures pour inverser cela s'il le souhaite.

Toutes ces mesures sont ancrées dans le désir de Cloudflare de donner du pouvoir aux créateurs de contenu. Le blog de Cloudflare souligne que les propriétaires de sites "devraient avoir le contrôle sur l'activité des bots IA sur leurs sites web" (Source: blog.cloudflare.com), et que robots.txt peut servir de "Code de Conduite" pour les bots (Source: blog.cloudflare.com). Mais parce que robots.txt seul repose sur le bon comportement, Cloudflare le complète par une application active (via son pare-feu) et des valeurs par défaut judicieuses. Comme l'a noté un analyste, le WAF (Web Application Firewall) de Cloudflare peut "appliquer ces règles" et bloquer les agents utilisateurs indésirables à la périphérie du réseau – une garantie bien plus solide qu'un simple fichier texte (Source: workmind.ai).

La démarche de Cloudflare offre ainsi à la fois **signal et application**. Les propriétaires de sites signalent "pas d'IA" via des robots et des paramètres mis à jour, tandis que le réseau périphérique mondial de Cloudflare peut effectivement refuser ou ralentir les crawlers non autorisés. Dans leur blog, Cloudflare se vante même que sa gestion des bots peut distinguer les crawlers humains des crawlers IA, appliquant les blocages en conséquence (Source: <u>adgully.me</u>).

En résumé, Cloudflare a élaboré une boîte à outils pour redonner le contrôle aux auteurs : des paramètres par défaut qui les protègent, ainsi que des options pour débloquer ou monétiser si désiré. Le raisonnement est succinctement énoncé par le PDG de Cloudflare : « Le contenu original est ce qui fait d'Internet l'une des plus grandes inventions », et il doit être « protégé » par un modèle économique qui fonctionne pour tous (Source: adgully.me).

Justification économique et éthique

Les principales justifications de Cloudflare pour le blocage par défaut des crawlers d'IA sont axées sur la **durabilité économique** et l'**équité numérique**. Les responsables soulignent à plusieurs reprises que l'ancienne économie du web, basée sur les clics, vacille sous le poids de l'IA. Comme l'a expliqué Matthew Prince, si les utilisateurs reçoivent des réponses des bots IA au lieu de cliquer, "l'incitation à créer du contenu original et de qualité [pour les sites] disparaît" et "l'avenir d'Internet est en péril" (Source: adgully.me). Le raisonnement est que les créateurs de contenu (journalistes, blogueurs, éducateurs) ont besoin de trafic pour monétiser leur travail. L'exploration par l'IA sans réciprocité menace cette source de revenus.

Les éditeurs eux-mêmes ont fait écho à cette logique. Par exemple, la News/Media Alliance (représentant plus de 2 200 éditeurs américains) a averti qu'ignorer les signaux "ne pas explorer" pourrait "saper la monétisation du contenu et l'industrie du journalisme" (Source: www.reuters.com). Des cadres supérieurs des médias comme Roger Lynch, PDG de Condé Nast, et Neil Vogel, PDG de Dotdash Meredith, ont salué la décision de Cloudflare, affirmant qu'elle créerait "un échange de valeur équitable sur Internet" et permettrait aux éditeurs de "limiter l'accès à notre contenu aux partenaires IA désireux de s'engager dans des accords équitables" (Source: adgully.me). Les grandes entreprises Internet — Reddit, Gannett, Pinterest, Ziff Davis — ont publiquement exprimé des points de vue similaires, présentant la politique de Cloudflare comme alignant les incitations à l'innovation et à la création de contenu (Source: adgully.me) (Source: <a href="adgul

Un autre aspect est l'**éthique des données** et l'idée de consentement. Le blog de Cloudflare et les commentaires associés soulignent que les utilisateurs ne réalisent souvent pas que leur contenu est collecté pour l'IA commerciale. Le blog de Workmind note que les propriétaires de sites "n'avaient aucune idée que leur travail acharné était utilisé pour construire des produits d'IA de plusieurs milliards de dollars" (Source: workmind.ai). La norme dominante — les bots peuvent collecter n'importe quoi à moins d'être explicitement bloqués — est contestée comme étant injuste. Beaucoup soutiennent que cela devrait devenir un scénario d'opt-in : les crawlers d'IA doivent respecter le consentement des créateurs (via robots.txt ou des contrats). La politique de Cloudflare impose ce changement.

Il y a aussi des implications légales. Bien que robots.txt ne soit pas lui-même juridiquement contraignant, Cloudflare souligne que les en-têtes dans les robots ou les chartes de licence pourraient acquérir une valeur juridique (Source: www.cloudflare.net). En rendant les signaux clairs et facilement accessibles, ils renforcent l'argument selon lequel les bots ont ignoré les préférences des propriétaires de sites à leurs propres risques. De plus, les poursuites intentées par de grands éditeurs (par exemple, NYT, AP, Rolling Stone) contre des entreprises d'IA soulignent que l'utilisation de données sans consentement relève de problèmes de droits d'auteur et de contrats (Source: apnews.com) (Source: www.reuters.com). L'approche de Cloudflare, qui exige une permission, peut aider à éviter de tels litiges en établissant un marché (ou un mécanisme de contrôle d'accès) autour du contenu web.

Enfin, il y a un argument d'**équilibre concurrentiel**. Cloudflare note que les entreprises d'IA (en particulier les grandes entreprises technologiques) peuvent simplement explorer le web sans frais, tandis que toute startup ou petit concurrent doit faire de même pour être compétitif. Le blocage par défaut "construit des clôtures" autour du web (selon les termes d'une analyse (Source: workmind.ai), forçant un nouvel équilibre. Ce faisant, la politique favorise sans doute un développement plus éthique de l'IA – encourageant les accords de licence et les partenariats de contenu plutôt que le parasitisme. En effet, l'initiative de Cloudflare encourage les développeurs d'IA à devenir des "partenaires" plutôt que des prédateurs sur le web ouvert (Source: adgully.me) (Source: workmind.ai).

En somme, le raisonnement de Cloudflare est que la viabilité à long terme du web exige de donner aux **propriétaires de contenu un réel choix et une compensation potentielle** pour l'utilisation des données. La politique de blocage par défaut est justifiée comme un correctif à un système asymétrique qui favorise actuellement les entreprises d'IA au détriment des créateurs.

Cas illustratifs et perspectives

Point de vue des éditeurs

Point de vue des éditeurs

Les grands éditeurs et les entreprises de médias numériques ont ouvertement soutenu les initiatives de Cloudflare. Par exemple, Condé Nast (éditeur de Vogue, Wired, etc.) a qualifié le blocage par défaut de « révolutionnaire » qui établit une nouvelle norme : les entreprises d'IA ne doivent plus s'approprier gratuitement le contenu (Source: <u>adgully.me</u>). La direction de USA Today Network a souligné qu'en tant que « plus grand éditeur du pays », le blocage du scraping non autorisé est « d'une importance capitale » pour protéger la propriété intellectuelle de valeur (Source: <u>adgully.me</u>). Ces voix considèrent la politique de Cloudflare comme une extension de leurs propres appels de longue date au respect et à la compensation.

Les organisations de licences applaudissent également ce changement. La déclaration de la Reuters News Media Alliance (Mt. [6]) a présenté le fait d'ignorer les robots comme une atteinte aux perspectives de monétisation. Le communiqué de presse de Cloudflare cite le PDG de l'Alliance qui vante l'outil de Cloudflare comme permettant aux éditeurs de toutes tailles de « reprendre le contrôle » de leur contenu (Source: www.cloudflare.net). De même, des agences comme le RSL Collective soutiennent que le contenu doit être non seulement protégé, mais aussi correctement sous licence et suivi, s'alignant sur les signaux techniques de Cloudflare (Source: www.cloudflare.net).

À un niveau plus granulaire, les petits créateurs de contenu et les professionnels du SEO ont noté des avantages techniques. Le scraping agressif par GPTBot et d'autres peut entraîner une augmentation de la charge du serveur et de l'utilisation de la bande passante. Le guide de Workmind souligne que le blocage de ces robots « protège les performances de votre site web » et réduit les coûts d'hébergement (Source: workmind.ai). De nombreux webmasters ont déjà activé l'interrupteur de blocage d'IA de Cloudflare pour cette raison (réduction des pics de charge) avant même de considérer les droits de contenu (Source: blog.cloudflare.com) (Source: blog.cloudflare.com).

Dans la jurisprudence, les éditeurs soulignent que l'entraînement d'une IA sans permission peut constituer une infraction. Par exemple, le scraping ouvert du web a conduit le New York Times à poursuivre OpenAl fin 2023 (Source: apnews.com). Le Times a soutenu que les réponses de ChatGPT (et la récupération « sans clic ») privaient le journal de revenus publicitaires et violaient ses droits d'auteur. La position de Cloudflare fait écho à cette lutte : elle offre aux propriétaires de sites un réglage par défaut intégré « sans scrapeurs », contournant l'ambiguïté juridique en empêchant l'action.

Point de vue des entreprises d'IA

Du point de vue des développeurs et chercheurs en IA, les changements de Cloudflare ont été controversés. Beaucoup dans le domaine de l'IA affirment que les modèles ont besoin de vastes données web et que l'exigence de permissions individuelles complique la collecte de données. Certains considèrent robots.txt comme un héritage qui ne devrait pas contraindre l'apprentissage automatique (surtout si les données sont accessibles au public). En effet, lorsque Cloudflare a accusé Perplexity d'ignorer robots.txt, l'équipe de Perplexity a vivement contesté, qualifiant cela d'argument de vente (Source: www.itpro.com). Ils soutiennent que le web a été conçu pour le crawling et que les robots devraient être libres d'accéder aux données publiques (invoquant souvent les doctrines de « fair use » dans les discussions juridiques) (Source: workmind.ai).

Les critiques soutiennent également que les mesures de Cloudflare pourraient « verrouiller » le contenu, entravant potentiellement l'innovation. Des commentateurs technologiques ont noté que l'exigence de paiements ou de permissions pourrait réduire la disponibilité des données pour des services d'IA bénéfiques (Source: www.techradar.com). Une analyse de TechRadar a averti que le système de paiement au crawl de Cloudflare « traite toutes les pages web comme ayant la même valeur » et pourrait décourager les entreprises d'IA, car d'énormes quantités de données web peuvent être obtenues à partir de sources publiques gratuites (comme Common Crawl) (Source: www.techradar.com). Si les entreprises d'IA sont confrontées à des coûts de licence complexes, les petites startups d'IA pourraient avoir du mal à collecter des données d'entraînement, ce qui renforcerait les acteurs établis ou les modèles soutenus par l'État. La critique est que « les systèmes actuels comme le paiement au crawl ne parviennent pas à résoudre le déséquilibre fondamental... la bataille sur les droits des données d'IA est plus une question de pouvoir que de paiement » (Source: www.techradar.com).

D'autre part, certains au sein de la communauté de l'IA reconnaissent que le passage à des modèles de permission est inévitable. Une vision équilibrée suggère que l'exigence d'accords ou de frais pour l'accès aux données pourrait professionnaliser les marchés de données. Dans le guide Workmind, la section « développeur d'IA » concède que, bien que les changements de Cloudflare compliquent la vie des créateurs d'IA, ils pourraient conduire à une IA plus éthique s'appuyant sur des sources de données bien documentées (Source: workmind.ai). De plus, l'industrie technologique dans son ensemble évolue vers des pratiques de données plus transparentes (par exemple, l'étiquetage de la provenance des données (Source: www.arrayanalytics.io), de sorte que la politique de Cloudflare pourrait accélérer l'établissement de normes.

En résumé, les entreprises d'IA présentent le contre-argument selon lequel des blocages généralisés pourraient étouffer l'innovation ou créer une disponibilité de données fragmentée. L'approche de Cloudflare impose un choix : soit se conformer aux propriétaires de sites, soit trouver des philosophies alternatives. L'affrontement avec Perplexity – au cours duquel Cloudflare a publiquement retiré le crawler de Perplexity de sa liste de « bots vérifiés » après la détection d'évasion (Source: www.itpro.com) – illustre cette tension. Il reste à voir comment les services d'IA s'adapteront (par exemple, en négociant l'accès, en développant des ensembles de données alternatifs ou en faisant pression pour des réglementations).

Point de vue des utilisateurs et services web

Du point de vue de l'utilisateur final, les effets sont subtils mais significatifs. À court terme, une conséquence de la politique de Cloudflare est que **l'ouverture du web est plus restreinte**. Les utilisateurs pourraient remarquer que certains futurs outils d'IA n'intègrent plus le contenu des sites qui refusent le crawling. Par exemple, si le contenu d'un site est bloqué, un outil de résumé d'IA pourrait ne plus répondre aux questions basées sur les articles de ce site. Pour les utilisateurs, cela pourrait signifier que certaines réponses deviennent moins complètes ou s'appuient sur moins de sources.

Cependant, de nombreux commentateurs de l'industrie s'attendent à peu de perturbations immédiates. Le guide Workmind note que les utilisateurs moyens « remarqueront un impact minimal » initialement (Source: workmind.ai) : le contenu n'apparaissant pas dans ChatGPT ou les nouvelles fonctionnalités de questions-réponses de Google ne nuit pas directement à un utilisateur, il ne fait que refuser des réponses basées sur l'IA à partir de ce contenu. Avec le temps, l'espoir est qu'une utilisation plus éthique des données améliore la confiance. Par exemple, si les entreprises d'IA doivent divulguer leurs sources ou payer pour du contenu de haute qualité, les utilisateurs pourraient en fait obtenir des réponses plus fiables et traçables à l'avenir.

Pour l'infrastructure web générale, cette politique met également en évidence une tendance vers un **web permissionné**. Les sites web exigent de plus en plus que tout crawler s'identifie et déclare ses intentions (recherche vs. analyse vs. entraînement). Cela pourrait conduire à des normes comme le protocole de permissions de Text and Data Mining (TDM) du W3C (Source:

Rankstudie

<u>www.arrayanalytics.io</u>), qui est conceptuellement aligné avec ce que fait Cloudflare. Pendant ce temps, Google (roi de la recherche) est sous pression pour séparer l'indexation de recherche traditionnelle de l'indexation d'IA – puisqu'il utilise « Googlebot » pour les deux (Source: <u>www.windowscentral.com</u>) (Source: <u>www.arrayanalytics.io</u>).

Globalement, si les clients de Cloudflare (propriétaires de sites) gagnent en contrôle, les fonctionnalités basées sur l'IA qui dépendent du crawling public pourraient devoir s'adapter. Les futures expériences de navigation ou de recherche pourraient évoluer : par exemple, si un utilisateur interroge un assistant IA, il pourrait recevoir des avertissements indiquant que certaines informations ne sont pas disponibles en raison de la protection du site. Comme l'a noté un analyste, l'écosystème dans son ensemble sera « meilleur lorsque le crawling sera plus transparent et contrôlé » (Source: <u>adgully.me</u>), ce qui pourrait bénéficier aux utilisateurs en clarifiant la provenance des informations.

Normalisation et Contexte Légal

Les actions de Cloudflare recoupent également des efforts plus larges visant à codifier les normes de crawling web. Plusieurs organismes de normalisation réagissent aux mêmes problèmes. L'IETF (Internet Engineering Task Force) est déjà en train de réviser le protocole robots.txt pour gérer les cas d'utilisation de l'IA (Source: www.arrayanalytics.io) (Source: <a href="www.arr

Le W3C (World Wide Web Consortium) a entrepris des travaux complémentaires. Son protocole de droits de Text and Data Mining (TDM) permet aux éditeurs de faire des déclarations lisibles par machine sur ce qui est autorisé en matière de fouille de données sur leur contenu (Source: www.arrayanalytics.io). Cela va au-delà de robots.txt en envisageant l'application technique des droits d'auteur ou des conditions de licence. La stratégie de Cloudflare fait écho à cela en rappelant aux entreprises la **signification juridique** des préférences de site (Source: www.cloudflare.net) (Source: www.arrayanalytics.io) - préparant essentiellement un avenir où les bots ne respectant pas robots.txt ou les règles TDM pourraient faire face à des réclamations contractuelles ou de droits d'auteur.

Sur le plan juridique, les régulateurs commencent tout juste à intervenir. Des décisions récentes (par exemple, le refus des régulateurs de données de l'UE d'arrêter l'entraînement de Llama de Meta sur les données d'Instagram (Source: www.infosecurity-magazine.com) montrent des résultats mitigés. Aux États-Unis, des affaires de droits d'auteur en cours (par exemple, Ziff Davis contre OpenAl (Source: www.reuters.com), Atlantic RM contre Microsoft) testent si le scraping de contenu publiquement disponible pour l'entraînement d'IA constitue un « fair use » ou une infraction. Les nouveaux signaux de Cloudflare, de par leur conception, créent des preuves de consentement ou de son absence (ce qui pourrait être important devant les tribunaux). Au minimum, l'entreprise estime que rendre les préférences explicites renforce les arguments de « rupture de contrat » contre les bots de scraping (Source: www.infosecurity-magazine.com) (Source: www.cloudflare.net).

Les critiques soutiennent qu'à moins que les législateurs n'agissent, les mesures purement techniques comme robots.txt n'ont pas de « mordant » exécutoire (même Cloudflare admet que ses politiques ne *garantissent* pas la conformité (Source: www.windowscentral.com). La discussion de l'IETF citée dans la <u>liste de diffusion</u> montre une certaine résistance à l'intégration de mandats exécutoires dans robots.txt, craignant que cela ne devienne une loi de facto (Source: mailarchive.ietf.org). Néanmoins, un changement à l'échelle de l'industrie (la règle par défaut de Cloudflare étant l'exemple principal) pourrait à lui seul créer une norme de facto. Déjà, des entreprises comme Microsoft (en partenariat avec Cloudflare sur des normes web « favorables à l'IA » (Source: www.techradar.com) et Google (avec des politiques de contenu similaires) se débattent avec la manière d'adapter leurs bots d'indexation.

En résumé, la politique de blocage par défaut de Cloudflare s'inscrit dans un paysage de gouvernance en évolution. Elle pourrait être complétée ultérieurement par des normes ou des lois formelles. Pour l'instant, l'application au niveau du réseau de Cloudflare est le mécanisme le plus immédiat pour concrétiser ce que les régulateurs et les organismes de normalisation commencent seulement à débattre.

Discussion: Implications et Orientations Futures

Implications Immédiates: La décision de Cloudflare modifie l'équilibre immédiat des pouvoirs sur le web. Les propriétaires de contenu sur le réseau de Cloudflare disposent désormais d'outils efficaces à portée de main. La majorité des sites hébergés dans le cloud peuvent rapidement se protéger contre le crawling d'IA indésirable. Les premiers indicateurs montrent que de nombreux propriétaires de sites ont déjà choisi volontairement de bloquer les bots d'IA (plus d'un million l'ont fait avec l'option de juillet 2024 (Source: adgully.me). Le nouveau réglage par défaut étend cette protection à pratiquement tous les nouveaux venus, évitant ainsi la nécessité de connaissances ou d'actions de la part de chaque propriétaire.

Pour les fournisseurs de services d'IA, l'implication est claire : ils doivent désormais demander l'accès aux comités. Certains pourraient interagir avec les sites via des API ou des accords de licence. D'autres pourraient se concentrer sur le contenu qui reste largement accessible. Nous pourrions assister à une prolifération de sites « favorables aux crawlers d'IA » qui optent volontairement (peut-être en échangeant des avantages contre de la visibilité) et de sites « résistants aux crawlers d'IA » qui protègent leur contenu. Le paysage pourrait se fragmenter.

Défis Potentiels : - Contournements de l'application : Des scrapeurs astucieux pourraient tenter de contourner les blocages de Cloudflare (par exemple, en faisant tourner les user-agents ou les adresses IP), tout comme certains essaient de contourner robots.txt aujourd'hui (Source: www.itpro.com). Cloudflare a renforcé la détection (retirant les contrevenants de sa liste de « bots vérifiés » (Source: www.itpro.com), mais des acteurs déterminés pourraient persister. Ce jeu du chat et de la souris suggère que le blocage par défaut pourrait n'être que partiellement efficace si les scrapeurs l'ignorent. Cependant, l'ampleur de Cloudflare (20 % du trafic web (Source: www.windowscentral.com) (Source: adgully.me) signifie que sa politique a toujours une portée étendue pour les acteurs conformes.

- Impact sur la recherche : La grande inconnue est la réaction des moteurs de recherche. Le double rôle de Google en tant que robot d'exploration de recherche et moteur de contenu IA complique les choses. Actuellement, un site ne peut pas différencier le « GoogleBot » utilisé pour le SEO du « GoogleBot » utilisé pour la collecte de données obscures (Source: www.windowscentral.com). Si de nombreux webmasters commencent à bloquer « GoogleBot » sans discernement pour protéger leur contenu, ils risquent de disparaître complètement de l'index de Google. Cloudflare reconnaît implicitement cette préoccupation; leurs recommandations suggèrent de bloquer Google-Extended (si distinct) plutôt que GoogleBot, mais cela est complexe et sujet aux erreurs (Source: www.xataka.com). Cette tension signifie que les propriétaires pourraient toujours être confrontés à un compromis entre visibilité et protection. La manière dont Google s'adaptera finalement (par exemple, en proposant des balises robots distinguant l'utilisation par l'IA) aura un impact considérable.
- Adoption des standards: Les signaux de contenu de Cloudflare dans robots.txt pourraient à terme gagner du terrain au-delà de la plateforme Cloudflare. L'entreprise a déjà mis en place une nouvelle « Politique de signaux de contenu » avec des balises spécialisées (ai-train, search, ai-input) et publie des outils pour encourager leur adoption (Source: www.cloudflare.net). Si l'IETF ou le W3C standardise des balises similaires, même les sites non-Cloudflare pourraient signaler aux robots d'exploration. Dans ce scénario, le blocage par défaut de Cloudflare deviendrait un exemple précoce d'une norme mondiale.

Perspectives à long terme : La grande question est de savoir si ces solutions technologiques seront suffisantes ou durables. Certains analystes sont sceptiques quant aux mécanismes comme le paiement par exploration, suggérant que des **stratégies juridiques et collectives** seront finalement nécessaires. La critique de TechRadar soutient que la monétisation seule ne résoudra pas le déséquilibre sans « levier » (action unifiée des éditeurs, lois applicables) (Source: www.techradar.com). En effet, certains éditeurs intentent des actions en justice en parallèle. Les outils de Cloudflare pourraient en partie servir de mesure provisoire pour démontrer la demande du marché, incitant les entreprises d'IA et les décideurs politiques à conclure des accords formels ou des réglementations.

À l'avenir, nous pouvons nous attendre à de nouvelles innovations. Cloudflare et ses partenaires explorent déjà l'authentification des agents (pour s'assurer que les robots d'exploration s'identifient honnêtement) et les licences structurées (par exemple, via le RSL Collective) qui automatisent les paiements ou exigent des rapports d'utilisation. Du côté des données, des technologies comme le suivi de la provenance du contenu (C2PA) pourraient compléter les règles d'exploration en filigranant l'origine du contenu. Si elles sont largement adoptées, celles-ci pourraient créer un écosystème où le contenu web ne peut pas être utilisé par les modèles d'IA sans attribution claire ou permission.

Cependant, certains experts s'inquiètent des effets secondaires. La restriction des robots d'exploration va-t-elle accélérer la nature « jardin clos » d'Internet ? Les chercheurs en open source et les universitaires trouveront-ils des sources de données alternatives, potentiellement moins réglementées ? La fragmentation pourrait-elle ralentir l'innovation ? L'interaction de ces forces se déroulera

sur plusieurs années.

Dans tous les cas, Cloudflare a affiché une position ferme : **les propriétaires de sites fixent les conditions d'engagement**. Comme l'a dit le PDG de Cloudflare, « les entreprises d'IA, les moteurs de recherche, les chercheurs et toute autre personne explorant des sites doivent être qui ils prétendent être. Et toute plateforme sur le web devrait avoir son mot à dire sur qui prend son contenu et pourquoi » (Source: <u>adgully.me</u>). Ce principe – transparence et consentement – est au cœur du changement de politique de Cloudflare.

Conclusion

La décision de Cloudflare de créer un robots.txt par défaut restreignant les robots d'exploration d'IA sur les nouveaux sites reflète un changement majeur dans la gouvernance du web, impulsé par l'IA générative. Leur raisonnement, fondé sur des données et amplifié par le plaidoyer des éditeurs, vise à **réaligner les incitations**: s'assurer que les créateurs continuent de bénéficier du trafic qu'ils génèrent, et exiger des systèmes d'IA qu'ils respectent la propriété du contenu. En passant d'un modèle d'opt-out à un modèle d'opt-in, Cloudflare place le contrôle explicite entre les mains des propriétaires de sites web.

Cette politique reconnaît que l'ancien modèle - « web ouvert signifie données d'entraînement librement disponibles » - est insoutenable pour un écosystème dynamique d'éditeurs indépendants. La suite d'outils de Cloudflare (boutons de blocage, robots.txt géré, signaux de contenu, paiement par exploration) constitue une stratégie holistique pour faire respecter cette nouvelle norme. Les premières données montrent un large soutien et une adoption par les éditeurs, tout en générant des réticences de la part de certains développeurs d'IA.

En substance, Cloudflare parie que le web ne peut pas survivre à l'ère de l'IA sans une **économie de contenu basée sur la permission**. Si cette position prévaut, nous pourrions voir un avenir où les données web sont traitées comme n'importe quelle autre ressource : sous licence et rémunérées. Alternativement, si le scraping incontrôlé se poursuit, le contenu des éditeurs pourrait simplement disparaître derrière des paywalls plus stricts ou des silos fragmentés.

Le résultat dépendra de nombreux facteurs : l'adaptabilité des entreprises d'IA, la réaction des moteurs de recherche, les décisions juridiques sur l'utilisation des données, et la manière dont la communauté web mondiale (sites sur et hors Cloudflare) réagira. Ce qui est clair, c'est que Cloudflare a jeté le gant. Leur blocage par défaut et leurs initiatives de robots gérés représentent un moment décisif – une note technique à un débat plus large sur les droits, l'utilisation équitable et l'avenir d'un internet ouvert.

Toutes les affirmations ci-dessus sont tirées de rapports sectoriels actuels, des publications de Cloudflare et de la couverture des événements en cours (Source: adgully.me) (Source: www.cloudflare.net) (Source: www.reuters.com) (Source: www.reuters.com). Ces sources documentent les données, les citations et les réactions qui soustendent les actions de Cloudflare et les arguments qui les entourent.

Étiquettes: cloudflare, robots-ia, robots-txt, extraction-donnees, protection-contenu, gptbot, grands-modeles-langage, economie-contenu

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.