AIボットと機関リポジトリ

COAR Annual Conference 2025での発表と議論

2025年9月8日 みんなでおさらいCOAR2025 第2回 「機関リポジトリ meets AI」

田辺浩介(物質・材料研究機構)

https://orcid.org/0000-0002-9986-7223

COAR Annual Conference 2025での「AIボットとリポジトリ」セッション

発表資料がCOARのWebサイトに掲載されています
 https://coar-repositories.org/news-updates/coar-annual-conference-2025/

Conference Presentations

Monday, May 12

Small group discussion about scaling multilingualism across the global repository network

Dealing with AI bots in repositories

Introduction to session and results of COAR survey - Kathleen Shearer, COAR

Current experiences of "friendly" network services - Petr Knoth, CORE Dynamic Bots Blocker, Lautaro Matas, LA Referencia

The Tale of IRD Friendly Robot, Paul Walk, COAR / Antleaf

Discussion: What approaches can repositories take to limit the harms of AI bots while remaining open to desirable value-added services?

AIサービスはどうやってWeb上の情報を収集しているのか

- Alサービスの運営者が、自動的にWebページのリンクをたどって情報を収集 するプログラムを実行している
- このプログラムを「ボット」「クローラー」と呼ばれる

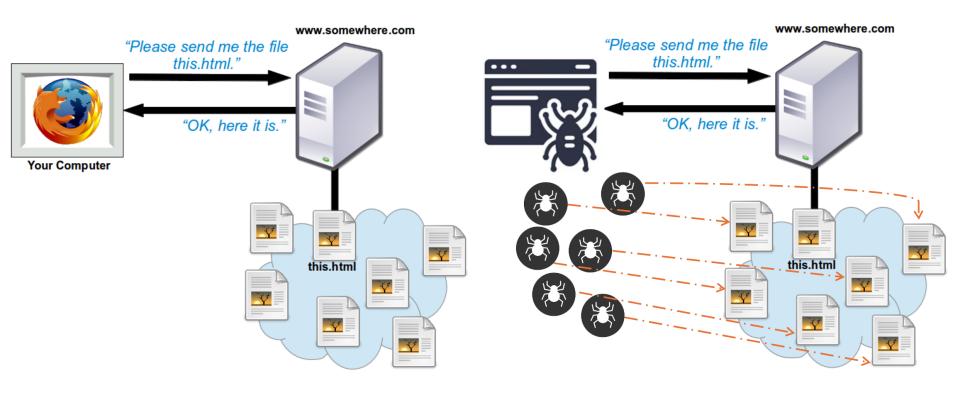
 検索エンジンのボットやクローラーと全く同じ
- この発表では、特にAIサービスの運営者が運用しているボットやクローラーを「AIボット」と呼びます

COARによる「AIボットとリポジトリ」に関連した調査・活動

- 2025年4月: AIボットによるリポジトリへの影響をたずねるアンケートを実施
- 5月: COAR Annual Conference 2025で"Dealing with AI bots in repositories"セッションを開催
- 6月:4月のアンケートの報告書を公開
- 7月: 「AIボットとリポジトリ」タスクフォース (Al Bots and Repositories Task Force)を設置
- 秋: タスクフォースによる報告書を公開予定

AIボットがなぜ問題視されるのか

- 人間よりもはるかに短期間で大量のアクセスを仕掛けてくる
 - 人間は興味のあるWebページのリンクだけをクリックするが、AIボットは片っ端 からリンク先をたどってWebページの中身をダウンロードする
- リポジトリのサーバーやインターネット回線に大きな負荷がかかり、 ときにサーバーが停止する
 - AIボットはメタデータの何十・何百倍ものサイズである論文や研究データ ファイルの中身を読むため、ファイル全体をダウンロードする必要がある

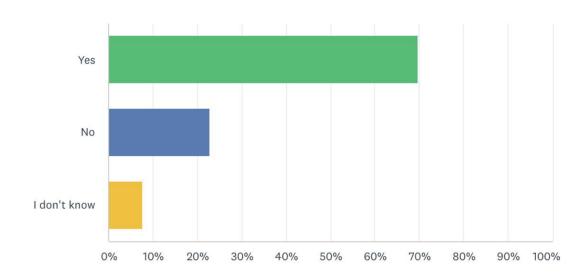


Ben Zhao. Dealing with Generative AI, Harms and Mitigation Techniques. Open Repositories 2025, Chicago, Illinois, 2025. https://doi.org/10.5281/zenodo.15790708

2025年4月のCOARのアンケートで、70%の回答者が 「AIボットによるリポジトリのサービス停止の経験あり」と回答

Have any AI bot encounters resulted in service disruption?

Answered: 66 Skipped: 0



Kathleen Shearer, Paul Walk. The impact of Al bots and crawlers on open repositories: Results of a COAR survey, April 2025. 2005. https://coar-repositories.org/wp-content/uploads/2025/06/Report-of-the-COAR-Survey-on-Al-Bots-June-2025-1.pdf

AIボットの運用者とコンテンツ提供者の非対称性

- AIボットの運用者は収集したコンテンツの内容で利益をあげるが、 コンテンツ提供者にはそれを還元しない
- コンテンツ提供者側にはボットによる不利益(サービス停止やその 対応の手間と費用)ばかりがかかる

困ってるならアクセスを拒否すれば?

ところが機関リポジトリにおいては、

アクセス拒否は簡単ではない

機関リポジトリのコンテンツ

- ほとんどの機関リポジトリのコンテンツは無料で提供されており、 コンテンツの対価を要求することが難しい
- 特に、多くのコンテンツがCreative Commonsライセンスで公開 されている

Creative Commonsライセンスとの関係

- コンテンツがCreative Commonsライセンスで公開されている場合、第三者(AIサービス業者など)に対してそのライセンスの範囲で自由に利用できることを許諾していることになる
 - 「AIサービスに使うならお金を払え」という利用条件を追加するのであれば、そのコンテンツはCreative Commonsでは提供できなくなる
- 現在のCCにはAI学習について制限をかける条項がないため、 人間だろうがAIだろうが(収集による)複製や再配布を拒めない
 - 拒めるのは商用利用(NC)や派生物の作成(ND)のみ
 - <u>CC Signals</u>という、AIに対してコンテンツの利用方法の意思表示を行うため の仕組みが開発されているが、まだ普及していない

じゃあCCで公開するのをやめる?

- おそらく現実的ではない
 - オープンアクセス・オープンサイエンスに逆行する動きを受け入れられるか?
 - AIサービスが研究成果を学習しなくなることを受け入れられるか?
- 日本の著作権法ではAIの学習のための複製(=ダウンロード)は著作者の 許諾なしで自由に行えるため、CC以外のライセンスで公開したとしても 著作権によってAIサービスの学習の制限をかけることはできない
- そもそも今回問題になっているのは「AIの学習」ではなく「AIボットの大量 アクセス」であり、仮にCCでの公開をやめて契約などで学習を禁じたところ で、AIボットの見境なしの大量アクセスを防ぐ手段にはならない

AIボットを制御する方法はないの?

- 検索エンジンのクローラー同様、Webサーバーに **robots.txt ファイル**を置く ことで、AIボットに収集頻度や対象の制限を「**お願い**」することは可能
- ただし**あくまで「お願い」でしかなく**、AIボットは**それを無視して**コンテンツにアクセスすることが可能であり、実際にそのような事例が多く見られる
 同様にAIボット向けのIlms.txtファイルがあるが、おそらくrobots.txtほど普及していない
- 自機関のリポジトリや学会・出版社のURLの".jp"や".com"などの後ろに "/robots.txt"をつけてアクセスしてみよう
 - 例: https://jpcoar.repo.nii.ac.jp/robots.txt

JAIRO Cloud@robots.txt



https://jpcoar.repo.nii.ac.jp/robots.txt

一般的なボット対策とその限界 (1)

- AIボットを名乗ってくるアクセスをブロックする
 - Webページにアクセスするソフトウェア (ブラウザーなど) は、アクセス時に ソフトウェアの名前 (Edge, Chrome, Firefoxなど) を名乗ることになっている
 - AIボットも名前を名乗ることになっている(GPTBot, ClaudeBotなど)
 - →AIボットがソフトウェアの名前を騙ってアクセスする例がある
- 短時間で多数のアクセスを行ってくるアクセス元をブロックする
 - IPアドレスを使用してブロック
 - →複数のIPアドレスから手分けしてアクセスしてきてしまう
 - →IPアドレスの拒否範囲を広げると、正当なアクセスを拒否してしまう
- アクセスしているのが人間であるかどうかの確認作業を入れる
 - 絵合わせパズルなどを解かせる→AIがパズルを解いてしまう

Webページにアクセスする際に名乗るソフトウェア名

● Firefoxの名乗るソフトウェア名(User-Agent)の例

```
Mozilla/5.0 (Windows NT 6.1; Win64; x64; rv:47.0) Gecko/20100101 Firefox/47.0 Mozilla/5.0 (Macintosh; Intel Mac OS X x.y; rv:42.0) Gecko/20100101 Firefox/42.0
```

ボットやクローラーの名乗るソフトウェア名の例

```
Mozilla/5.0 (compatible; Googlebot/2.1; +http://www.google.com/bot.html)
```

```
Mozilla/5.0 (compatible; YandexAccessibilityBot/3.0; +http://yandex.com/bots)
```

User-Agent. https://developer.mozilla.org/ja/docs/Web/HTTP/Reference/Headers/User-Agent

ScienceDirectのrobots.txt (一部)



https://www.sciencedirect.com/robots.txt

AIサービス業者の取り組み

- AIボットで使用するソフトウェア名やIPアドレスを開示している
- しかし、これでAIボットのアクセスを除外できるかというと、そうでもない

GPTBot

GPTBot is used to make our generative AI foundation models more useful and safe. It is used to crawl content that may be used in training our generative AI foundation models. Disallowing GPTBot indicates a site's content should not be used in training generative AI foundation models.

OpenAIのボット が名乗る名前

Full user-agent string: Mozilla/5.0 AppleWebKit/537.36 (KHTML, like Gecko); compatible; GPTBot/1.1; +https://openai.com/gptbot

Published IP addresses: https://openai.com/gptbot.json

Overview of OpenAl Crawlers. https://platform.openai.com/docs/bots

AIボットがソフトウェア名を偽っているとされる例

Our multiple test domains explicitly prohibited all automated access by specifying in robots.txt and had specific WAF rules that blocked crawling from Perplexity's public crawlers. We observed that Perplexity uses pot only their doclared user agent but robots.txtで指定されていないソフトウェア名 also a generic browser intended to imperso their declared crawler was blocked. 「お願い」を無視したことにはならない

Declared	Mozilla/5.0 AppleWebKit/537.36 (KHTML, like Gecko; compatible; Perplexity-User/1.0; +https://perplexity.ai/perplexity-user) 本当の名前		20-25m daily requests
Stealth			3-6m daily requests

Gabriel Corral, Vaibhav Singhal, Brian Mitchell, Reid Tatoris.Perplexity is using stealth, undeclared crawlers to evade website no-crawl directives. The Cloudflare blog. 2025.

https://blog.cloudflare.com/perplexity-is-using-stealth-undeclared-crawlers-to-evade-website-no-crawl-directives/

一般的なボット対策とその限界 (2)

- 悪質なAIボットの拒否リストを作る
 - COUNTER(電子ジャーナルの利用統計の規格)では、アクセス数の統計からボットのアクセスを除外するため、ボットのリストを作成している
 - NIIでもJAIRO Cloud向けにJAIRO Crawler-Listを作成している
 - →そのリストを誰がどうやってメンテナンスするのか?COUNTERのボットリストですら最終更新は2024年4月、JAIRO Crawler-Listに至っては最終更新は2015年12月
 - 前述のように、ボットは嘘の名前(ふつうのWebブラウザーなど)でアクセス してくることがあるため、リストだけではボットのアクセスを拒否しきれない



機関リポジトリのAIボット対策による副作用

- COARで開発している<u>IRD</u> (International Repositories Directory、 世界のリポジトリー覧)では、その作成のためにボットを使っている
- しかし、IRDのボットによる機関リポジトリへのOAI-PMHへのアクセスが 拒否されてしまう例が確認されている
- OAI-PMHはコンテンツ収集を行うボット(クローラー)のために用意されているのに、誤った(過剰な)ボット対策でそれが機能しなくなっている

Checking OAI-PMH

- No HTTP response (2653)
- 403 Not Authorised (111)
- 401 Unauthorized (8)

Webブラウザではアクセスできるのに、 ボットでアクセスすると"Not Authorized" エラーを返すリポジトリが複数存在する

Paul Walk. The tale of IRD, the friendly robot..... COAR Annual Conference 2025. 2025. https://coar-repositories.org/wp-content/uploads/2025/05/4.-The-tale-of-IRD-friendly-robot-Walk.pdf

人間向けとAI向けで提供方法を変える

- 通常の画面ではAIボットを弾くかわりに、AIボットにはサービスに影響しない別の方法でコンテンツを提供する
 - <u>WikipediaはKaggleとコンテンツ提供の契約を交わしている</u>
 - プレプリントサーバーのarXivは、<u>ダンプファイル(収録している論文の全ファイルと全メタ</u> <u>データ)を利用者が手数料を払う形で提供している</u>
 - ただし、これも名前を詐称するAIボットの対策にはならない
- 現在のところ、最も効果的なのは商用のWebアプリケーションファイアウォールを使ってボットのアクセスを弾くこと
 - 当然ながら、いいお値段がする

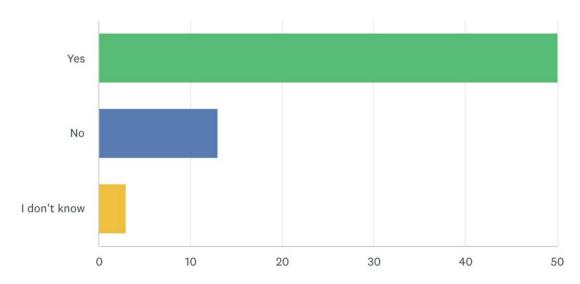
本気で対策しようとすると...

- 前述の複数の対策を組み合わせることになる
- 商用のWebアプリケーションファイアウォールは、それらを上手に 組み合わせてAIボットのアクセスを弾いてくれる
 - <u>世界中のボット一覧や「認証済みボット」のデータベース</u>を用意していることもある
- もちろん完璧ではないし、そもそも商品なのでコストがかかる

回答者の76%が「何らかのAIボットのアクセスの拒否・制限をしている」と回答

Have you applied any measures to minimize or stop access to the repository by AI bots?

Answered: 66 Skipped: 0



Kathleen Shearer, Paul Walk. The impact of Al bots and crawlers on open repositories: Results of a COAR survey, April 2025. 2005. https://coar-repositories.org/wp-content/uploads/2025/06/Report-of-the-COAR-Survey-on-Al-Bots-June-2025-1.pdf

「既定でAIボットをブロック」するファイアウォール

公平な取引の代わりに、現在のWebはAlクローラーによって搾取されており、コンテンツクリエイターにはほとんどトラフィックも価値ももたらされていません。

それが、本日7月1日、「コンテンツ独立記念日」と私たちが呼ぶ日をもって変わります。Cloudflareは、世界の主要な出版社やAI企業の大多数とともに、AIクローラーがクリエイターに報酬を支払わない限り、デフォルトでブロックする方針へと移行します。そのコンテンツこそがAIエンジンを動かす燃料であり、したがってコンテンツクリエイターが正当に報酬を受け取るのは当然のことです。

商用ファイアウォール提供企業のブログ。AIボット対策としては ありがたいが、オープンアクセスとの整合性を取る必要がある

Matthew Prince. コンテンツ独立記念日:報酬なしのAIクロールは許さない!. Cloudflareブログ. 2025. https://blog.cloudflare.com/ja-jp/content-independence-day-no-ai-crawl-without-compensation/

JAIRO CloudでのAIボット対策

- NIIサービス説明会2023の資料に 「BOT対応」が明記されている
- 「WAF、robots.txtによる対策」とあるが、 前述のとおりrobots.txtはほとんど空
- WAF (Webアプリケーションファイアウォール) でなんらかの対応をしているのかもしれないが、詳細は不明

JAIRO Cloud

NII RCOS

情報基盤センターの皆様に向けて

独自ドメイン利用機関を除き、 基本的に図書館単独運用を想定したサービス

- 問合せ対応: 1次切り分け後、通常問合せはJPCOARが対応
 - 窓口支援業者が1次切り分け、通常問合せはJPCOARが対応。緊急時はNII主導で対応を実施
- アクセス制御:リポジトリ管理者が機能、コンテンツへのアクセス制御
 - ユーザ認証:統合認証利用時は機関のアカウント管理はリボジトリ管理者が実施。
 学認IdP経由(試験提供)の場合は権限管理をリボジトリ管理者が実施
 - アクセス制御:アイテム、インデックス単位での公開・非公開制御、時限付き非公開設定、IPア
- セキュリティ対策
 - 運用監視およびWeb Application Firewall (WAF)による対応
- BOT対応
 - WAF、robots.txtによる対策
 - リポジトリ機能としてのアクセスログ記録時にBOT判定を実施。集計の際にBOTアクセス除外

スの他

• クラウドチェックリストの導入にむけて対応開始

国立情報学研究所. 公開基盤の新機能紹介. NIIサービス説明会2023. 2023. https://www.nii.ac.ip/openforum/upload/5-

3_setsumeikai2023_rcos_koukaikiban20231026.pdf

COAR Annual Conferenceでの議論

- リポジトリのコンテンツは、AIボットに対してもオープンであるべき
- 一方で、悪質なAIボットに対する防御策は必要
 - アクセスログをもとに、AIボットと思われるアクセスを一時的に拒否するプログラムが紹介 されていた
 - https://github.com/lareferencia/lareferencia-bot-blocker
- この点について、コミュニティとしてできることはないか?
 - 田辺の意見: 悪質なボットやIPアドレスのリストの共有くらいしか思い浮かばないが、 COUNTERのボットリストでもあまり更新されていないことを思うと、まともに機能しない 気がする。少なくとも、人手が介在する仕組みでは絶対無理だろう

まとめ

● AIボットとリポジトリの問題は、

「オープンアクセス・オープンサイエンスを支えるコストを、 誰がどのように工面するか」

という、巨大でまだ解決策のない世界的な課題のひとつ