

図書館ウェブサイト の公開性

クローラに対するアクセス
制御に関する調査

安形輝(亜細亜大学)

agata@asia-u.ac.jp

お知らせ

岡崎市立中央図書館

お知らせ

- 2010.9.7 週刊誌の貸出期間の変更について
- 2010.9.6 『子ども図書室だより』最新号(2010年9月号)を発行しました。
- 2010.9.1 岡崎市立中央図書館のホームページへの大量アクセスによる障害について
- 2010.8.18 テーマ展示のお知らせ

- お知らせ
- 図書館について
- 図書館カレンダー
- 蔵書検索
- 資料情報
- 地域資料
- デジタルアーカイブ
- マイページ
- こどもとじょしつ
- ティーンズコーナー
- 関連リンク
- サイトマップ



English | 中文 | 한국어 | Português |

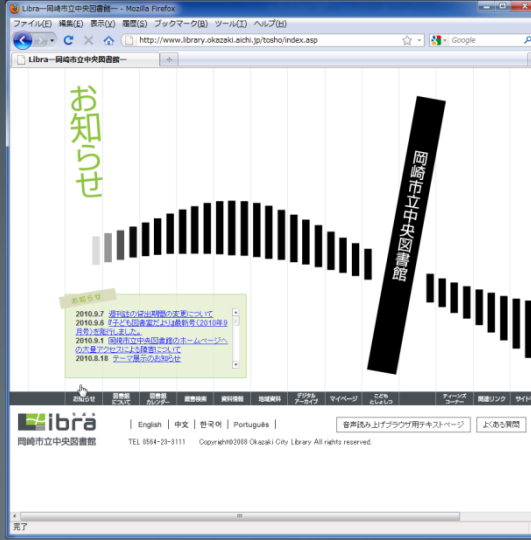
音声読み上げブラウザ用テキストページ

よくある質問

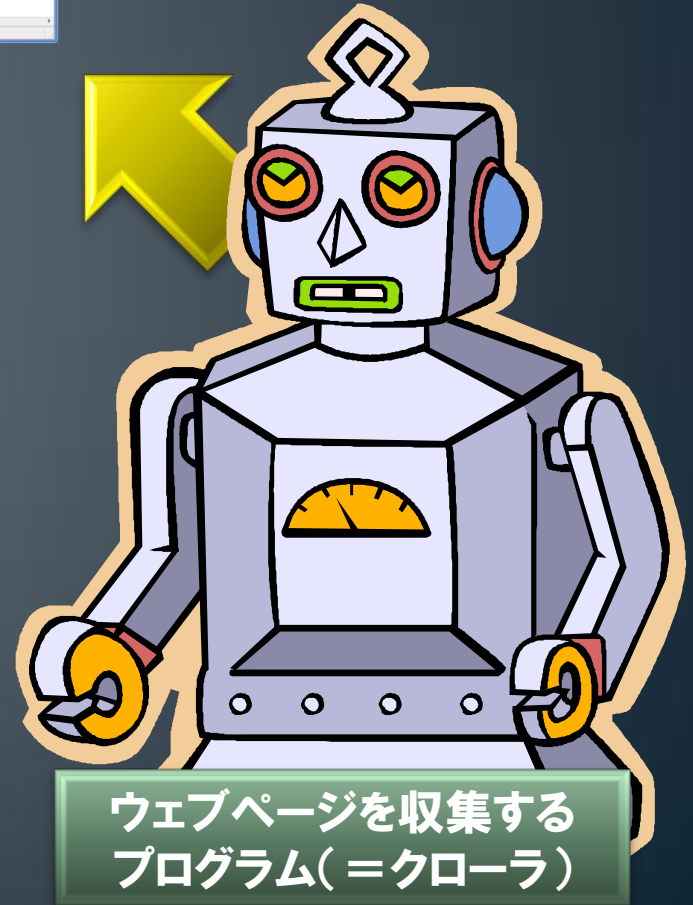
TEL 0564-23-3111 Copyright©2008 Okazaki City Library All rights reserved.

重要な広報手段

サービスアクセス手段



利用者



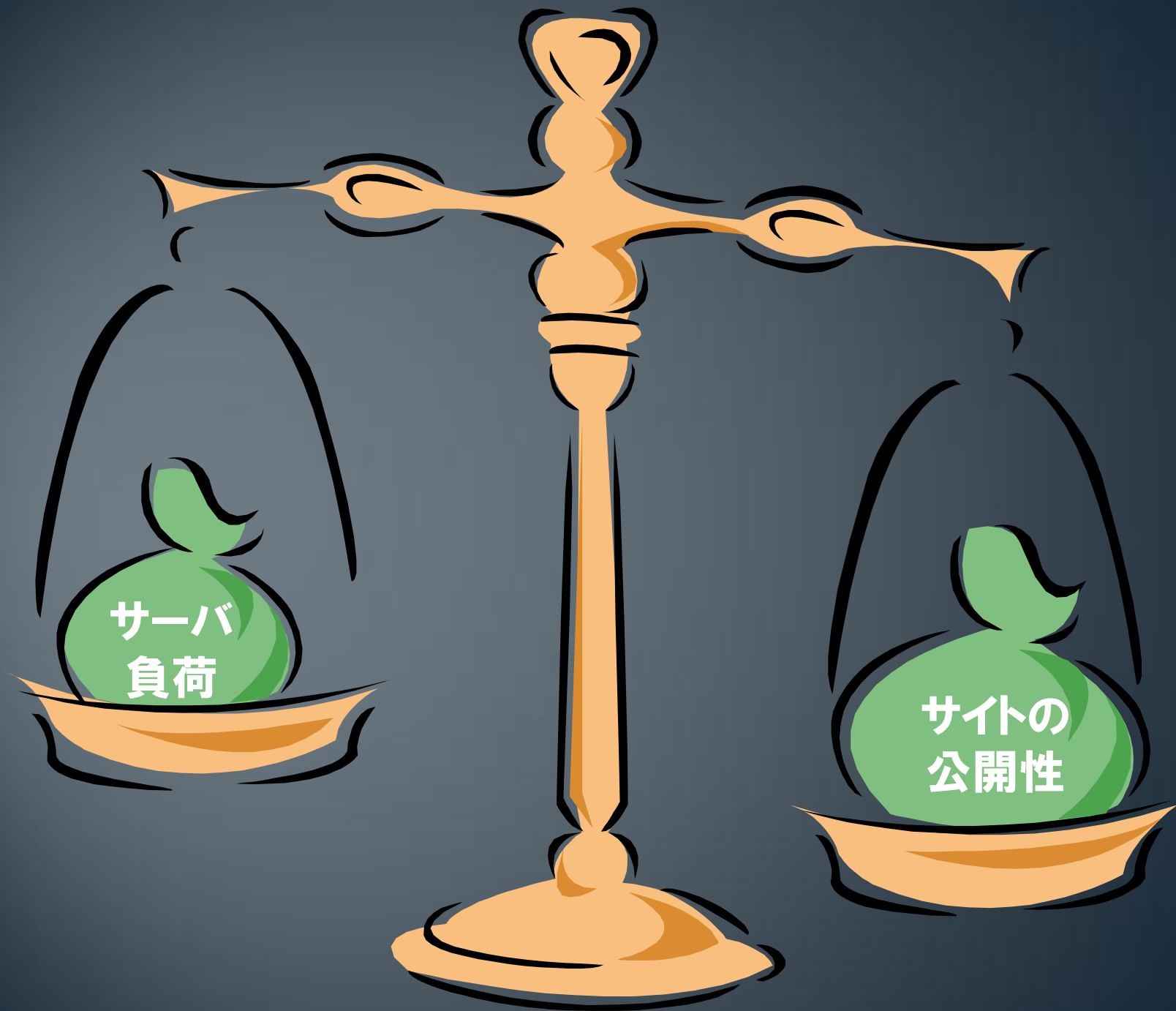
ウェブページを収集する
プログラム(=クローラ)



2010年6月



りぶら
ibra
岡崎市図書館交流プラザ



サーバ
負荷

サイトの
公開性

宮田洋輔 (慶應義塾大学)†
池内淳 (筑波大学)
†miyayo@slis.keio.ac.jp

安形輝 (亜細亜大学)
上田修一 (慶應義塾大学)

A Large-Scale Study of Robots.txt

Yang Sun, Ziming Zhuang, and C. Lee Giles
The Pennsylvania State University
University Park, PA, USA
{ysun, zzhuang, giles}@ist.psu.edu

抄録: 機関リポジトリに収録された文献の少なくない数が、深層ウェブ化していることが明らかになっている。そこで本研究では、その原因を明らかにするために、日本の機関リポジトリとリポジトリに収録された学術情報のアクセス可能性に関する調査をおこなった。本調査の結果から、robots.txt によって、検索エンジンからのアクセスを排除している事例があること。また、pdf ファイルのテキスト抽出の可否、全文 URL の長さなどの要因が、学術情報へのアクセスの可能性を低めていることが示唆された。

ABSTRACT

Search engines largely rely on Web robots to collect information from the Web. Due to the unregulated open-access nature of the Web, robot activities are extremely diverse. Such crawling activities can be regulated from the server side by deploying the Robots Exclusion Protocol (robots.txt). Although it is not all robots (and many confined in robots.txt). With our study of 7,593 websites from educational and business domains. Five crawls to study the temporal changes of the data, we present a summary of robots.txt at the Web scale. The use of robots.txt has increased

detail, especially at the scale of the Web. A study of usage of robots.txt in UK universities and colleges identified 163 websites and 53 robots.txt [2]. Robots.txt were examined in terms of file size and the use of Robots Exclusion Protocol within the UK university domains. We studied the usage of robots.txt as an aid for indexing examples from Fortune Global 500. We present the first large-scale study of robots.txt in the domains of education, government, and business. We present our observations and compare our data with previous studies.

INTRODUCTION

We collected the initial URLs to form the Open Directory Project (DMOZ). Our study covered three domains: education, government, and business. The university domain is further divided into the American, European, and Asian university domains. We use the Fortune Top 1000 Company List as our data source in the business domain. Our crawl performed five crawls for the same set of websites between Dec. 2005 and Oct. 2006.

3. RESULTS

Statistics: We crawled and investigated 7,593 websites including 600 government websites, 2,047 newspaper websites, 1,487 USA university websites, 1,420 European university websites, 1,039 Asian university websites, and 1,000 company websites.

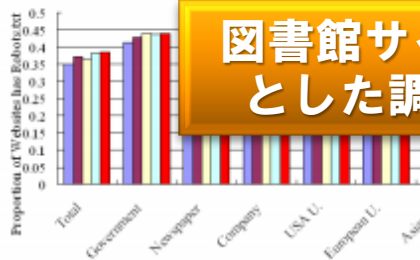


Figure 1: Probability of a website that robots.txt in each domain.

インターネット全体が対象のクロウラのアクセス排除に関する調査は多い

機関リポジトリについては、我々の研究グループで行った例がある

図書館サイトを対象とした調査はない

1. はじめに

機関リポジトリは大学や研究所による学術情報資源の公開・蓄積のために設置されている。Lynch は、機関リポジトリを「機関とそのコミュニティの構成員によって作成された電子資源の管理と発信のために、大学がそのコミュニティの構成員に提供する一連のサービス」と定義し、「広く一般の人々に向けてそのサービスを提供する」と述べている。

現在、日本では、国立情報学研究所の JAIRO が構築・運用されている。

機関リポジトリに収録された学術情報へのアクセス手段として、1)学術情報に直接、2)機関リポジトリ経由、3)検索エンジン経由、4)横断検索システムなど

学術情報資源に直接アクセスする場合、利用者は、何らかの形で情報資源のウェブ上での識別子 URL を知っている必要がある。

機関リポジトリ経由でアクセスする場合、利用者はそのリポジトリの存在を知っており、URL をアドレスバーに直接入力するか、ブックマークなどによってアクセスする。

機関リポジトリの学術情報にアクセスする場合は日常的なウェブ上で検索エンジンでキーワードを入力し、検索結果のなかで機関リポジトリに収録された情報資源と遭遇する可能性がある。その場合には、リポジトリ内の情報が、各種の検索エンジンによって、登録されている必要がある。

横断検索システムなどによるアクセスの場合、機関リポジトリが、たとえば OPEN DOAR³⁾や

OAister⁴⁾、日本の国立情報学研究所の JAIRO のようなサービス・プロバイダーに登録し、メタデータのハーベスティング(刈り取り)を可能な状態にしておく必要がある。

佐藤らによる機関リポジトリのアクセスログ研究によると、機関リポジトリに収録された文献の半数近くは検索エンジンを経由でおこなわれている。また、収録された学術情報の存在しているにもかかわらず検索エンジンで検索できないウェブ)化も指摘された。佐藤らは、元データを用いて検索した McCown⁵⁾ では、機関リポジトリのメタ

データ中の半数程度しか、検索エンジンによって、カバーされていないことが明らかになっている。また安形ら⁶⁾がおこなった、日本の機関リポジトリのメタデータから抽出した全文ファイルの登録状況の調査でも、Google, Yahoo!, Bing のいずれかの検索エンジンに登録されたものが7割、佐藤らの調査⁵⁾でもっともアクセスの割合が多かった Google だけでは53.2%と、検索エンジン経由でのアクセスが十分に機能していないことが明らかになっている。

そこで、本研究では、機関リポジトリに収録された学術情報へのアクセスの問題の要因を明らかにするために、機関リポジトリ自体とそこに収録された学術情報ファイルの調査をおこなった。

2. 調査の概要

表1に、1,000件以上の全文データを持った機関リポジトリでの、Google, Yahoo!, Bing のいずれかの検索エンジンからのアクセス可能な割合の上位10機関を、表2に下位10機関を示した。表から、検索エンジンでの登録率が100%のリポジトリも存在するものの、登録率が高いリポジトリでも必ずしも100%

図書館サイトは クローラの

アクセスを
排除してい
るか

その影響は
あるか

図書館サイトはクローラの
アクセスを排除しているか

ロボット排除プロトコル

複数の方法 / robots.txtが事実上の標準

例えば、robots.txtの記述例

全コンテンツ
排除

```
User-agent: *  
Disallow: /
```

一部コンテンツ
排除

```
User-agent: *  
Disallow: /cgi-bin
```

調査対象館

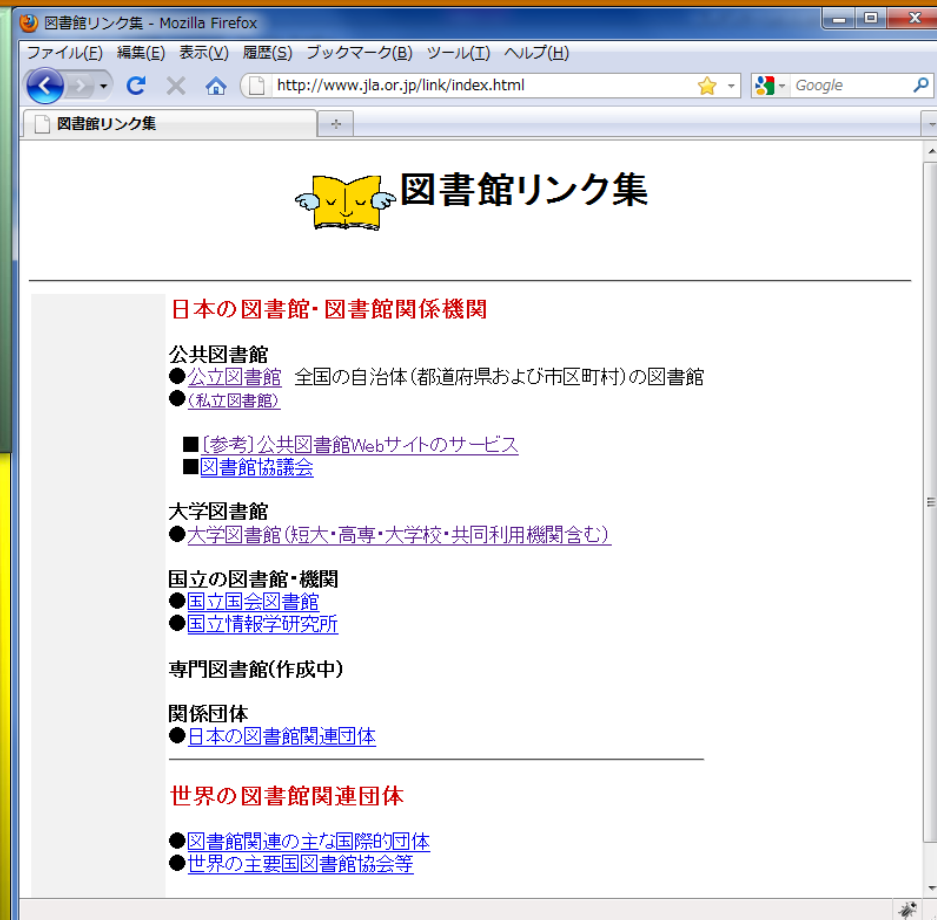
日本図書館協会図書館リンク集

<http://www.jla.or.jp/link/>

公共図書館
合計2,450館
大学図書館

2010年9月5日にアクセス可能だった2,065館

公共図書館	1,277館
大学図書館	788館



The screenshot shows a web browser window titled "図書館リンク集 - Mozilla Firefox". The address bar displays "http://www.jla.or.jp/link/index.html". The page content includes a logo of a yellow book with wings and the text "図書館リンク集". Below the logo, there are several sections of text:

- 日本の図書館・図書館関係機関**
 - 公共図書館
 - [公立図書館](#) 全国の自治体(都道府県および市区町村)の図書館
 - [\(私立図書館\)](#)
 - [参考] [公共図書館Webサイトのサービス](#)
 - [図書館協議会](#)
 - 大学図書館
 - [大学図書館](#) (短大・高専・大学校・共同利用機関含む)
 - 国立の図書館・機関
 - [国立国会図書館](#)
 - [国立情報学研究所](#)
 - 専門図書館(作成中)
 - 関係団体
 - [日本の図書館関連団体](#)
- 世界の図書館関連団体**
 - [図書館関連の主な国際的団体](#)
 - [世界の主要国図書館協会等](#)

調査の観点

robots.txt
の有無

- ・ ない場合、クローラのアクセス排除をしていない

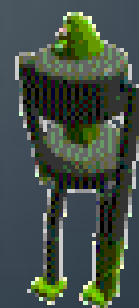
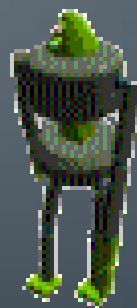
⇒ クローラOK

robots.txt
の内容

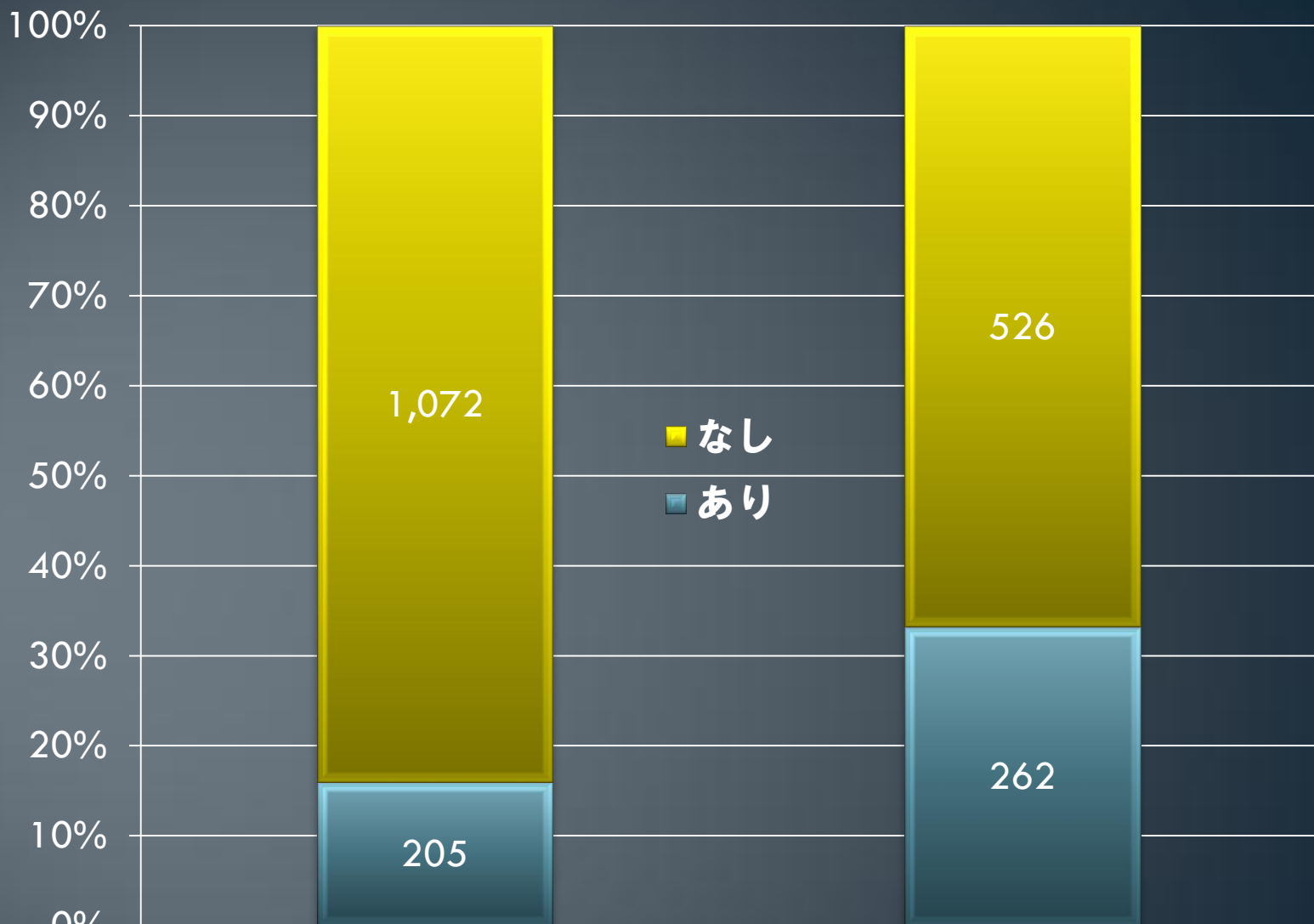
- ・ 誤りがあるか

アクセス制
限の対象

- ・ クローラのアクセスを全コンテンツで排除しているか



Robots.txtの有無



	公共図書館	大学図書館
■ なし	1,072	526
■ あり	205	262

誤りがあるか

Robots.txtの内容



ウェブページを返す



記述の誤り

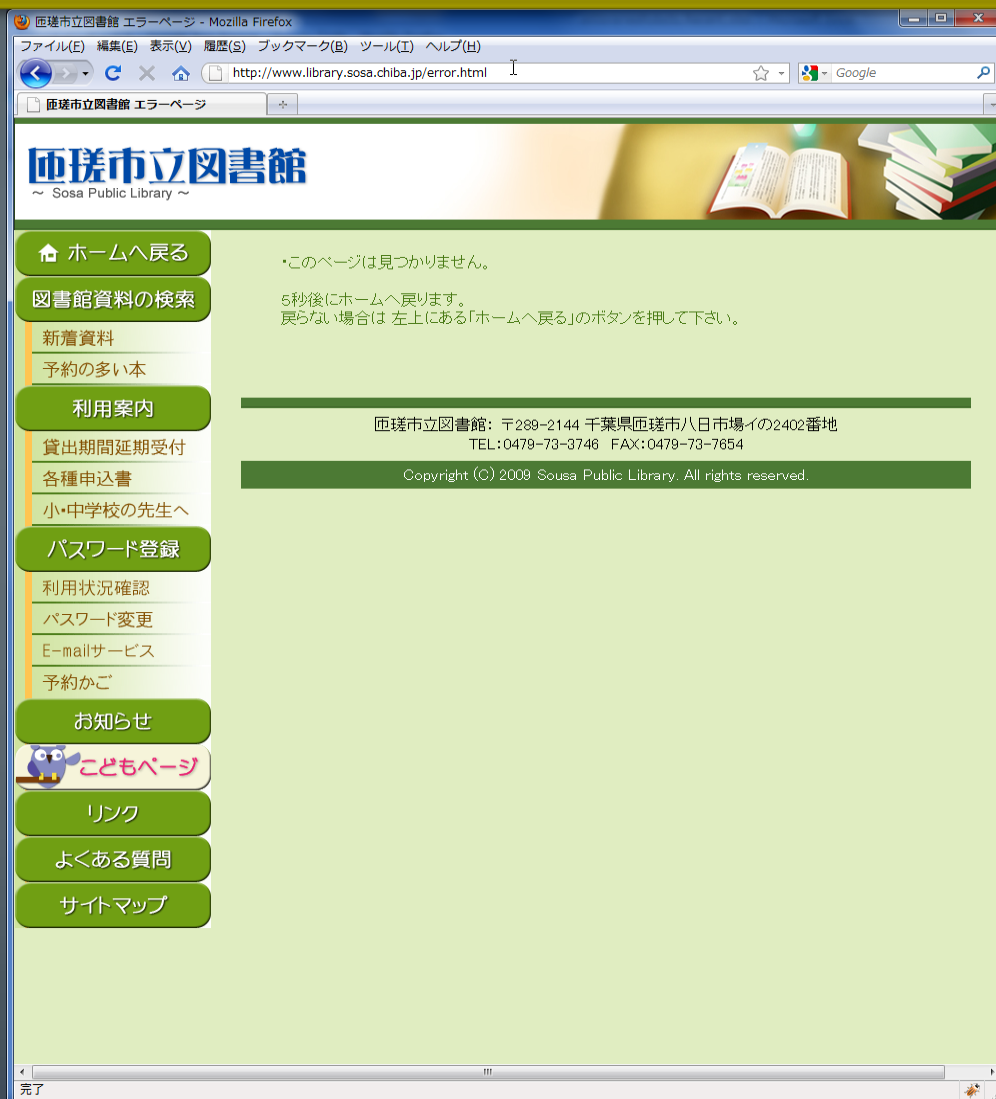
ウェブページを返す

Robots.txtの内容

robots.txtをリクエストするとウェブページを返す

ないのと同じ

⇒クローラOK 



記述の誤り



岡崎市立図書館を始めとする17館



誤り(記述に矛盾)で全アクセス排除



同一の業者によるサイト



誤りが修正されている館が5館

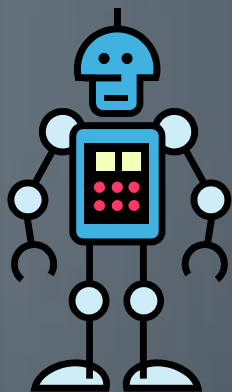
アクセス制限の対象

	公共図書館		大学図書館		全体	
	館数	割合	館数	割合	館数	割合
全公開・大半公開	1,216	95.2%	778	98.7%	1,994	96.6%
全排除	61	4.8%	10	1.3%	71	3.4%
合計	1,277	100.0%	788	100.0%	2,065	100.0%

図書館はクローラのアクセス制御をしているか

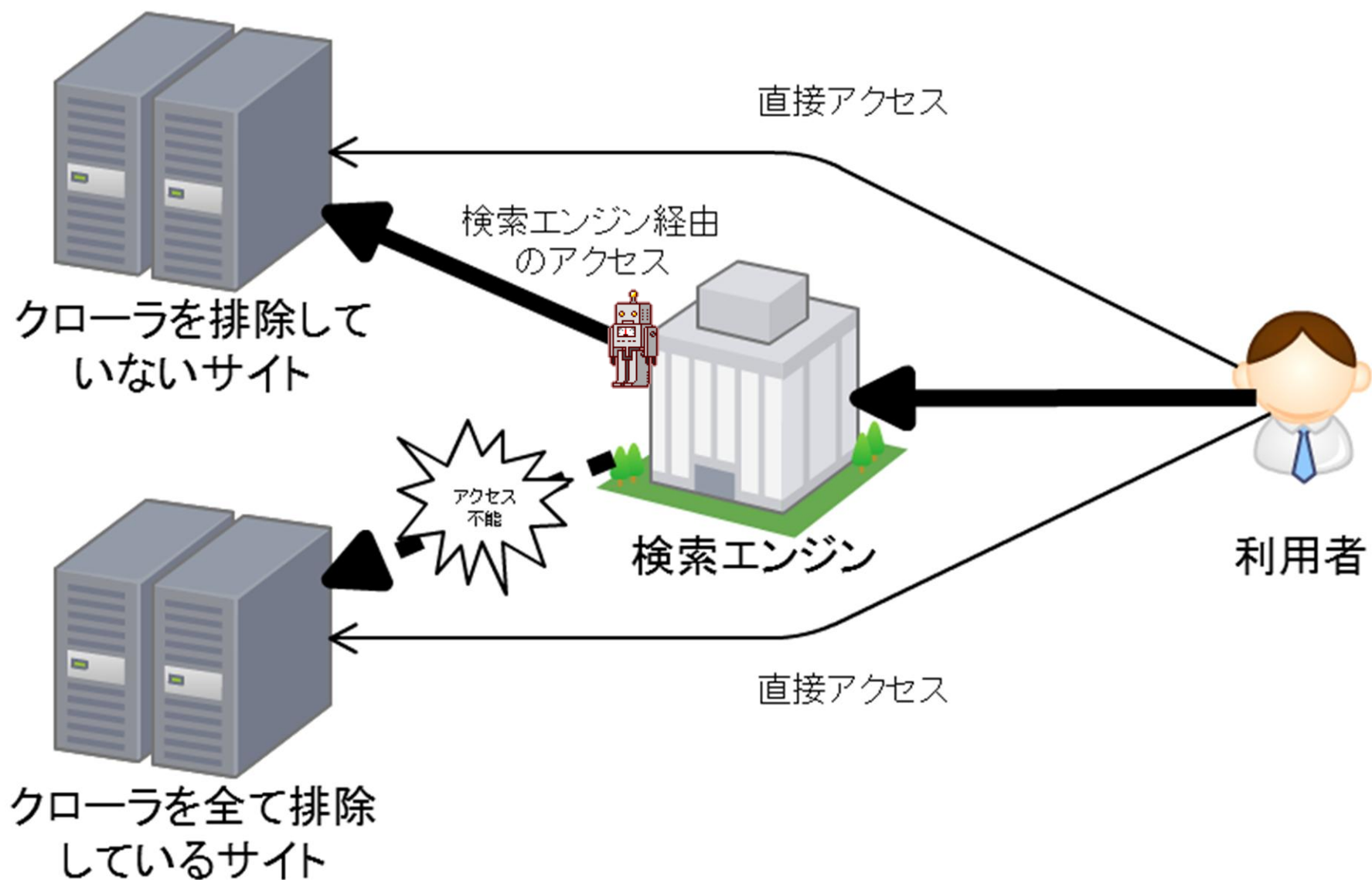
クローラのアクセスを

ほとんどの図書館
は全て・大半の部
分で認めている



一部の館は
全排除

クローラへのアクセス排除の影響



調査方法

検索対象



YAHOO!
JAPAN



Google
日本

検索式

図書館の
名称

検索結果

タイトル、
URL、要約

上位8位
まで

ox

表示(V) 履歴(S) ブックマーク(B) ツール(I) ヘルプ(H)

http://www.google.co.jp/webhp?hl=ja

「いなべ市図書館」の検索結果 - ... x +

ニュース 書籍 Gmail その他 ▾

Google 日本

いなべ市図書館

Google 検索 I'm Feeling Lucky

検索オプション
言語ツール

広告掲載 Google について Google.com in English

© 2010 - プライバシー

ox

履歴(S) ブックマーク(B) ツール(I) ヘルプ(H)

http://www.yahoo.co.jp/

Yahoo! JAPAN

オークション My Yahoo! YAHOO! JAPAN ツールバー ショッピング まっず

ウェブ 画像 動画 ブログ 辞書 知恵袋 地図 一覧 ▾

いなべ市図書館 検索

次世代ラーメン決定戦、投票受付中 ▶ 大沢たかお、田中圭らが語る「デキる男と女」 ▶ エベレスト登頂へ、栗城史多を応援しよう

トピックス 経済 エンタメ スポーツ その他

15時27分更新

- 仙谷氏 日中会談見送り認める
- 米車業界「為替操作」と批判
- 架空増資 柴野元議員を逮捕へ **NEW!**
- 270人271脚 正直にギネス断念
- ポル・ポト派元幹部4人を起訴
- イチロー無安打「M11」のまま
- 親が選ぶイクメン、つるの1位
- 田代容疑者、質問に震え出す

今日の話題(34件) 一覧

マリンスの恩返し
9月16日9時の分配信
神奈川新聞社

今日の天気
--% | --℃/--℃
表示する地域を指定 ▾

今日の予定 カレンダーを活用

今日の運勢 牡羊座 ▾ 71点 3 4 5

友だちをつくればもっと楽しめる!

賢くためる、使う

18日 11時35分

マクドナルド
アカギや「聖騎士星矢」が見
簡単レシピ264品でおいし

クローラのアクセスを全て排除している図書館をGoogleとYahooで検索すると・・・

特集 赤い彗星・シャア専用ザクがガンダムゾーンに出現 一覧

全国のうまいラーメン店といえば
ラーメン店ランキングを毎日発表。次世代
ラーメンの投票も!



郵便番号または市区町村名を入力してください。
指定した地域周辺の情報が表示されます。
例:「1060032」「港区」「六本木駅」など

【ニュース】表参道で「岡崎」伝統産業展(みん経)

「いなべ市図書館」の検索結果 - Mozilla Firefox

編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(I) ヘルプ(H)

http://www.google.co.jp/search?q=いなべ市図書館&ie=utf-8&oe=utf-8&q=t&rls=org.mozilla:☆

いなべ市図書館 - Google 検索

「いなべ市図書館」の検索結果 - ... x +

動画 地図 ニュース 書籍 Gmail その他 ▾

Google

いなべ市図書館

約 23,600 件 (0.03 秒) 検索オプション

[図書館TOPへ戻る - いなべ市](#)

[www.city.inabe.mie.jp/book/](#) - 類似ページ

[いなべ市図書館 蔵書検索](#)

詳しくはお持ちの携帯電話の取り扱い説明書をご覧ください。QRコード未対応の場合は、下記URLを入力して接続してください。http://lib.city.inabe.mie.jp/liswing/we/opaci/kensaku.jsp. QRコード: いなべ市TOPへ戻る・図書館TOPへ戻る.

[lib.city.inabe.mie.jp/liswing/we/opac/index.html](#) - キャッシュ - 類似ページ

[三重県の図書館](#)

伊勢市立小俣図書館, 519-0505, 伊勢市小俣町本町2. いなべ市員弁図書館, 511-0202, いなべ市員弁町楚原940. いなべ市大安図書館, 511-0274, いなべ市大安町大井田1305. いなべ市藤原図書館, 511-0511, いなべ市藤原町市場493-1. いなべ市北勢図書館 ...

[www.reference-net.jp/lib_dir/24.html](#) - キャッシュ - 類似ページ

[いなべ市大安図書館\(いなべ市図書館\)【町コミガイド】](#)

いなべ市大安図書館(いなべ市図書館). 住所: 三重県いなべ市大安町大井田1305. 電話番号: 0594-87-0021. アクセス(電車): 「大安駅」より. 駐車場: あり. 営業時間(期間): 9:30~17:30. 休業日: ...

[mie.town.co.jp/guide/as001694/](#) - キャッシュ - 類似ページ

[いなべ市 図書館:マピオン電話帳モバイル](#)

三重県いなべ市にある図書館の電話番号・地図・住所などを検索できるマピオン電話帳.

[p.mapion.co.jp/M13006/24214/](#) - キャッシュ

[いなべ市図書館](#)

いなべ市図書館 蔵書検索等システム操作方法. ●本の検索をするには 検索条件を入力し、「検索

「いなべ市図書館」の検索結果 - Yahoo!検索 - Mozilla Firefox

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(I) ヘルプ(H)

http://search.yahoo.co.jp/search?p=いなべ市図書館&search.x=1&fr=top_ga1_sa&tid

いなべ市図書館 - Google 検索

「いなべ市図書館」の検索結果 - ... x +

Yahoo! JAPAN - Yahoo!検索

ウェブ | 画像 | 動画 | ブログ | 辞書 | 知恵袋 | 地図 | 一覧 ▾

いなべ市図書館 [条件を指定して検索](#)
[検索設定](#)

ウェブ検索結果 いなべ市図書

[本を探します](#)

いなべ図書館蔵書検索. 操作方法. 1.メールアドレス変更. 2.パスワード変更
[lib.city.inabe.mie.jp/liswing/we/opaci/kensaku.jsp](#) - キャッシュ

[本を探します](#)

図書 雑誌. 対象館. 全館. 指定した館のみ. 北勢 員弁 大安 藤原. 一覧件数. 10 25 50 100. 論理式. 検索項目. キーワード. 一致条件. 著者名. 単独条件(半角) 絞り込み条件(半角)【新着資料一覧】 ...

[lib.city.inabe.mie.jp/liswing/we/opac/kensaku.jsp](#) - キャッシュ

[いなべ市役所ホームページ](#)

いなべ市のプロフィール、観光マップ等。
[www.city.inabe.mie.jp](#) - ブックマーク: 12人が登録

[いなべ市 図書館 :マピオン電話帳](#)

全国各地の主要スポット900万件の電話番号・地図・住所などを検索できるマピオン電話帳。ここは三重県いなべ市 図書館 のカテゴリです。三重県いなべ市にある図書館の詳細情報を地図や一覧から選んでチェック!

[www.mapion.co.jp/phonebook/M13006/24214](#) - キャッシュ

[いなべ市役所藤原図書館 - いなべ市 - 0594-46-4150 ...](#)

いなべ市役所藤原図書館. 住所. 〒511-0511 三重県いなべ市藤原町市場493-1 ... いなべ市役所藤原図書館について ... いなべ市役所藤原図書館以外のお店に、順番に電話をかけて行ったり、地図等の詳細情報を見に行く事が出来ます。 ...

[www.lococom.jp/ft/22430420825](#) - キャッシュ

[三重県の図書館](#)

伊勢市立小俣図書館, 519-0505, 伊勢市小俣町本町2. いなべ市員弁図書館, 511-0202, いなべ市員弁町楚原940. いなべ市大安図書館, 511-0274, いなべ市大安町大井田1305. いなべ市藤原図書館, 511-0511, いなべ市藤原町市場493-1. いなべ市北勢図書館 ...

[いなべ市の図書館 - MAPPLE 地図「ちず丸」](#)

いなべ市の図書館一覧から簡単に地図検索、ドラッグスクロール・範囲拡大・フリースケールで直感地図操作のMAPPLE 地図「ちず丸」 ... 三重県いなべ市大安町大井田1305. いなべ市藤原図書館. 住所: 三重県いなべ市藤原町市場493-1. いなべ市 ...

[www.chizumaru.com/czm/objlist-24214G0120.htm](#) - キャッシュ

完了

Googleでは一番最初に検索されるが、
Yahooでは10位以内には検索されない



クローラによるアクセス

排除なし

全排除

結果数

割合

結果数

割合

1

1601

80.3%

0

0.0%

2

167

8.4%

1

1.4%

3

34

1.7%

0

0.0%

4

11

0.6%

0

0.0%

5

2

0.1%

0

0.0%

6

5

0.3%

0

0.0%

7

2

0.1%

0

0.0%

8

0

0.0%

0

0.0%

検索結果順位

9位以下

172

8.6%

70

98.6%

合計

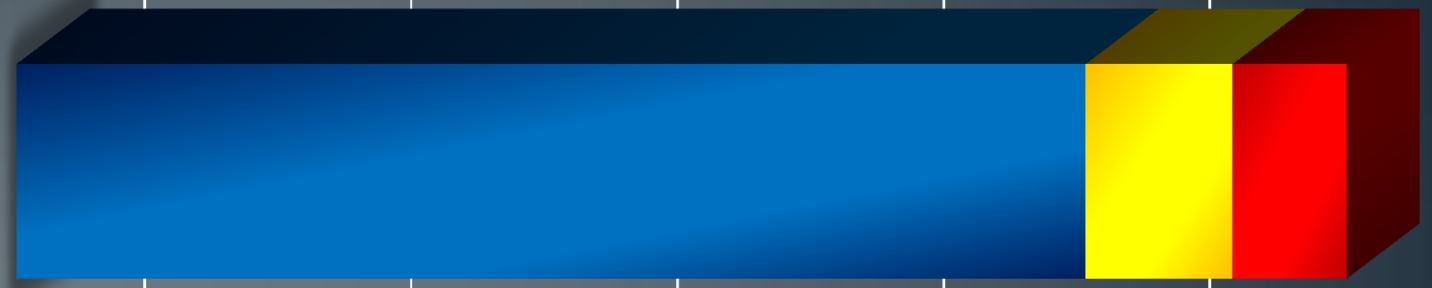
1994

100.0%

71

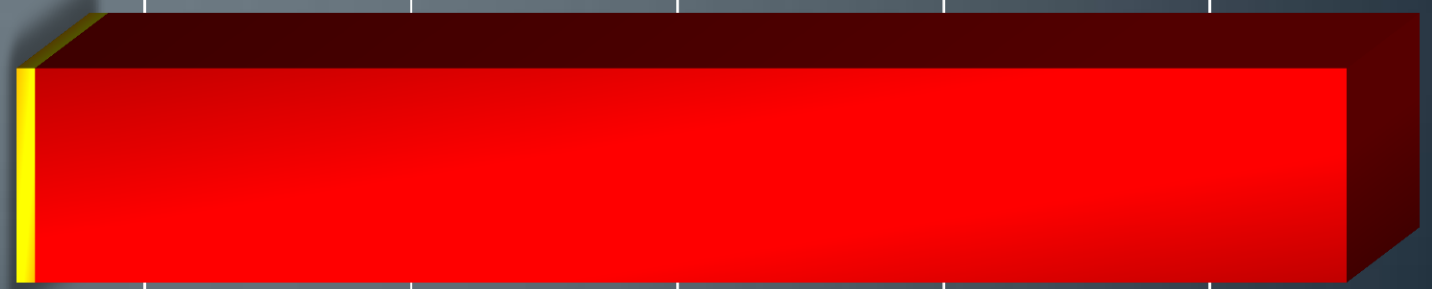
100.0%

排除なし



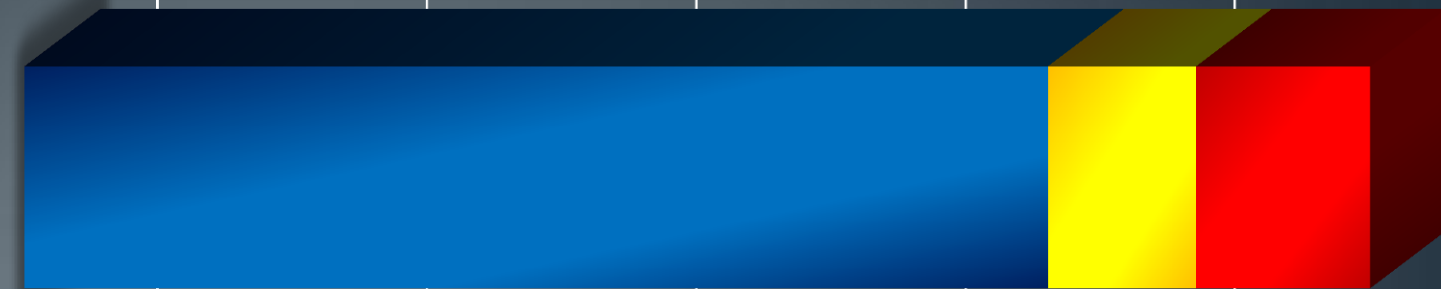
■ 1位 ■ 8位以内 ■ 9位以下あるいは検索されず

全排除



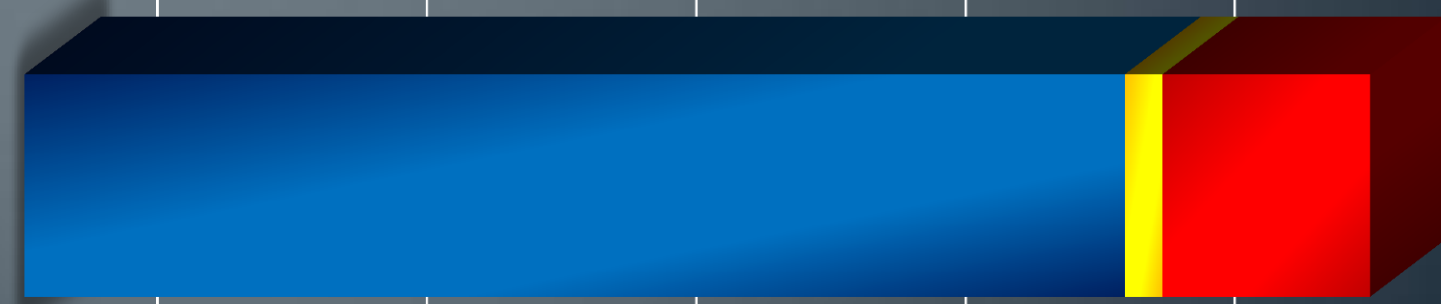
0% 20% 40% 60% 80% 100%

排除なし



■ 1位 ■ 8位以内 ■ 9位以下あるいは検索されず

全排除



0% 20% 40% 60% 80% 100%

全排除のページをなぜ Googleは検索できる？



Libra (りぶら) 岡崎市図書館交流プラザ

休館日: 祭日・休館日 一律休館
開館時間: 8:30~21:00
駐車料金: 延滞料金は60分間まで無料
*(窓口)に駐車券を提示して下さい。)

サイトマップ フロアマップ 周辺マップ
施設利用案内 図書館 生涯学習情報 施設予約システム 交通のご案内

金額	夜間	全日	延長時間
	8:00 ~ 9:00	8:30 ~ 13:00	22
	2:00 ~ 22:00	9:00 ~ 17:00	1時
			18:00
300円	16,520円	1,010円	1,800円
0,700円	21,290円	1,300円	2,320円
90円	1,040円	70円	110円

リンク情報を利用

図書館の配列は 公共図書館(愛知県、岐阜県、静岡県)の順です。
各県の大学図書館内の配列
→ 大学図書館、大学共同利用機関の配列

受付館名	担当窓口	TEL	FAX	e-mail	住所	受付方法	受付時間	図書部HP
コレクショ...	サービス課	052-212-5...	052-212-3...	at@libra.jp	〒490-0001 岡崎市古町	MAIL・FAX・郵送(書式類)・電話	10:00~20:00(平日) 10:00~18:00(土曜・日曜)	http://www.aichi-pref-library.jp/

愛知県図書館 WEB関連

【対応できるフレンスの範囲等】
【依頼時の留意点】
※依頼で購置した出典(資料)等の
#R84545

図書部HP <http://www.aichi-pref-library.jp/>
CPAC <http://www.aichi-pref-library.jp/cgi-bin/Snspress/shin/0/mode=1>

港区立図書館

約 261,000 件 (0.10 秒)

港区立図書館 サイト・トップページ

港区立図書館・国民読書年・お問い合わせ・English. 最終更新日: 2010年9月7日. みなと図書館・三田図書館・赤坂図書館・高輪図書館・港南図書館・麻布図書サービスセンター・郷土資料館. みなと図書館 ...

www.lib.city.minato.tokyo.jp/ - キャッシュ - 類似ページ

資料を探す	赤坂図書館
みなと図書館	高輪図書館
三田図書館	港南図書館
所在地・休館日・開館時間	図書館の利用案内

minato.tokyo.jp からの検索結果 >

港区立図書館 | みなと図書館

港区立図書館・国民読書年・お問い合わせ・English・図書館TOP・みなと図書館・三田図書館・赤坂図書館・高輪図書館・港南図書館・麻布図書サービスセンター・郷土資料館・図書館TOP > みなと図書館 ...

www.lib.city.minato.tokyo.jp/f/minato.cgi - キャッシュ - 類似ページ

【東京図書館制覇!】港区立図書館

港区立図書館の一覧、地図、データをまとめています。... 東京タワーの北東にある、港区立図書館の中央館。地域資料が充実しています。地域資料室への通路にある展示コーナーでは、図書館が所蔵する貴重な資料を使った展示が行われています。...

tokyo-toshokan.net/00000024.htm - キャッシュ

港区立図書館

港区の図書館へすると、blogs.y

港区立図書館

全国各地の主要スポット900万件の電話番号・地図・住所などを検索できるマピオン電話帳。ここでは東京都港区 図書館 のカテゴリです。東京都港区にある図書館の詳細情報を地図や一覧から選んでチェック!

www.mapion.co.jp/...> 公共施設, 図書館, 東京都 - キャッシュ - 類似ページ

港区立図書館でも月曜開館! : 日本創新党 荒川区議会議員小坂英二の ...

タイトル

要約

いなべ市図書館

約 23,600 件 (0.03 秒)

いなべ市図書館 TOPへ戻る - いなべ市

www.city.inabe.mie.jp/book/ - 類似ページ

いなべ市図書館 蔵書検索

詳しくはお持ちの携帯電話の取り扱い説明書等をご覧ください。QRコード未対応の場合は、下記URLを入力して接続してください。http://lib.city.inabe.mie.jp/iliswing/we/opaci/kensaku.jsp. QRコード. いなべ市TOPへ戻る・図書館TOPへ戻る.

lib.city.inabe.mie.jp/iliswing/we/opac/index.html - キャッシュ - 類似ページ

三重県の図書館

伊勢市立小俣図書館, 519-0505, 伊勢市小俣町本町2. いなべ市員弁図書館, 511-0202, いなべ市員弁町楚原940. いなべ市大安図書館, 511-0274, いなべ市大安町大井田1305. いなべ市藤原図書館, 511-0511, いなべ市藤原町市場493-1. いなべ市北勢図書館 ...

www.reference-net.jp/lib_dir/24.html - キャッシュ - 類似ページ

いなべ市大安図書館(いなべ市図書館)【町コミガイド】

いなべ市大安図書館(いなべ市図書館). 住所: 三重県いなべ市大安町大井田1305. 電話番号: 0594-87-0021. アクセス(電車): 「大安駅」より. 駐車場: あり. 営業時間(期間): 9:30~17:30. 休業日: ...

mie.town.co.jp/guide/as001694/ - キャッシュ - 類似ページ

いなべ市 図書館:マピオン電話帳モバイル

三重県いなべ市にある図書館の電話番号・地図・住所などを検索できるマピオン電話帳.

p.mapion.co.jp/M13006/24214/ - キャッシュ

いなべ市, いなべ市員弁図書館, 検索可, 有 ...

https://idx.milai.pref.mie.jp/.../EntryState.jsp - キャッシュ - 類似ページ

いなべ市図書館

いなべ市図書館 蔵書検索等システム操作方法. ●本の検索をするには 検索条件を入力し、「検索

排除していない港区立図書館の結果と全てを排除している図書館の結果

図書館サイトなのに危険な サイトに似ている？

正式名称
でないタ
イトル

要約が
ない



Google検索結果中の要約

クローラによるアクセス

排除なし

全排除

結果数

割合

結果数

割合

要約あり

1,505

99.3%

11

19.0%

要約なし

10

0.7%

47

81.0%

計

1,515

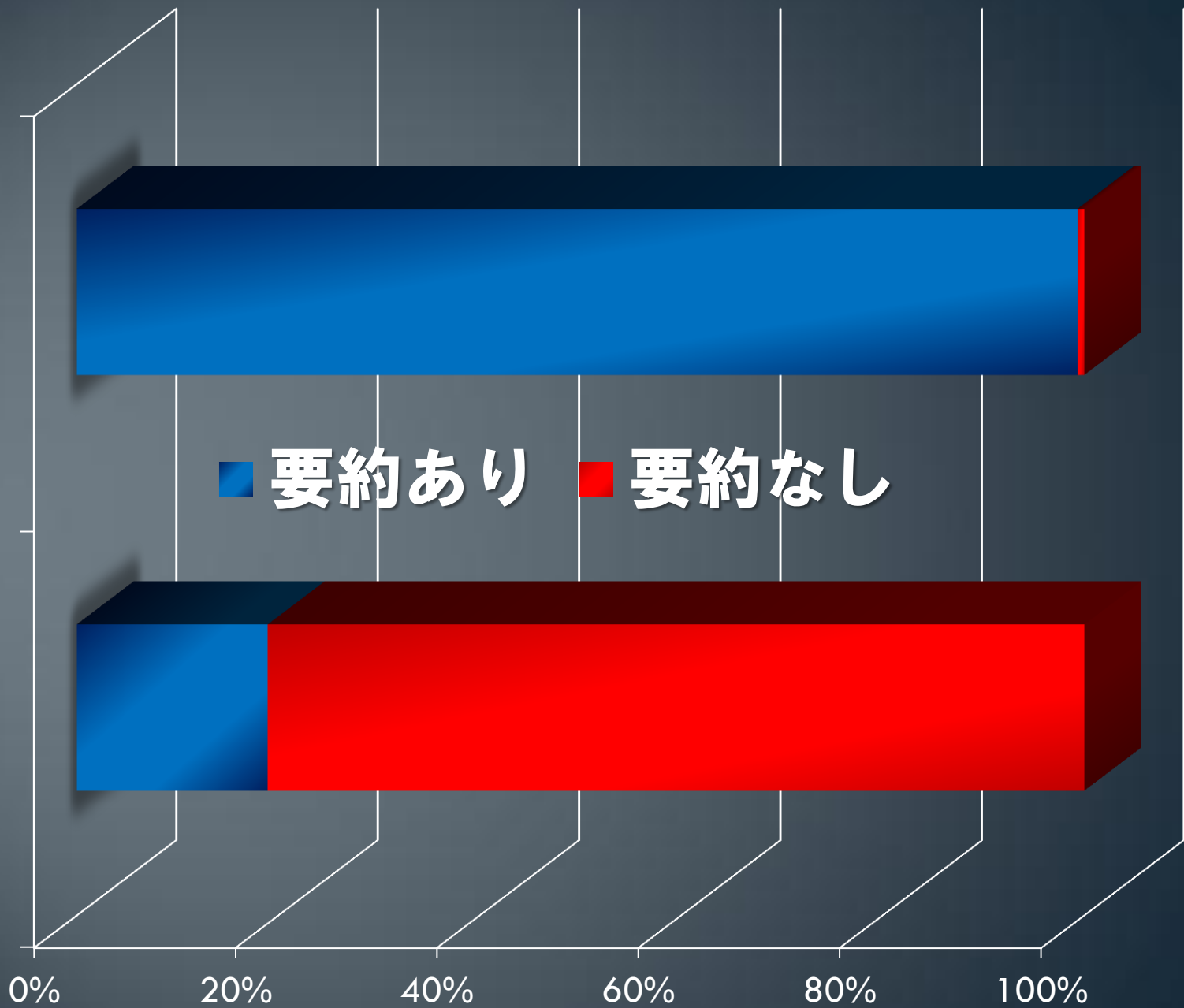
100.0%

58

100.0%

排除なし

全排除



クローラを排除したときの影響はあるか

検索結果に表示され
なくなる

The logo for Yahoo! Japan, featuring the word "YAHOO!" in a large, red, stylized font with a registered trademark symbol, and the word "JAPAN" in a smaller, red, sans-serif font below it.

表示されても危険な
サイト風になる

The logo for Google Japan, featuring the word "Google" in its multi-colored font, with the Japanese characters "日本" (Japan) in a smaller, black font below it.



まとめ



一部を除き多くの館はクローラを排除していない

全コンテンツで排除している館は検索しにくい状態に

図書館ウェブサイト の公開性

クローラに対するアクセス
制御に関する調査

安形輝(亜細亜大学)

agata@asia-u.ac.jp