

WWWコンテンツ統計調査報告書

～我が国の Web 上のコンテンツ情報量から見たインターネットの発展～

平成 16 年 7 月

総務省 情報通信政策研究所

主任研究官 佐伯 千種
研究官 島田 博也
研究官 田畑 伸哉

目次

はじめに - Web の役割と Web コンテンツ情報量調査の意義	4
1 . 調査の概要	6
1 -(1) 調査の仕組み・原理	6
1 -(2) 調査対象	8
2 . 調査結果について	9
2 -(1) これまでの調査実施状況	9
2 -(2) 調査結果	10
2 -(2)-1 2004 年 2 月 (第 11 回) 調査結果	10
2 -(2)-1-1 J P ドメイン総コンテンツ量	10
2 -(2)-1-2 J P ドメイン総ファイル数のファイル種別シェア	11
2 -(2)-1-3 J P ドメイン総データ量のファイル種別シェア	12
2 -(2)-2 過去の調査結果との比較	13
2 -(2)-2-1 J P ドメイン総サーバ数の推移	13
2 -(2)-2-2 J P ドメイン総ページ数の推移	15
2 -(2)-2-3 J P ドメイン総ファイル数の推移	16
2 -(2)-2-4 J P ドメイン総データ量の推移	18
2 -(2)-2-5 J P ドメイン 1 万ページあたり平均データ量と J P ドメイン 1 サーバあたり平均ページ数の比較	20
2 -(2)-2-6 J P ドメイン総コンテンツ量発展推移 (年別)	21
3 . 最近の調査結果から見られた新たな現象	22
3 -(1) 調査設計に関する変化	22
3 -(1)-1 J P ドメインの国内シェアの変化	22
3 -(1)-2 ロボットで取得できないマルチメディアデータの増加	23
3 -(1)-3 Web 規模の増大と調査インフラ増強の問題	23
3 -(1)-4 ネットワークのトラフィック増大によるハブの渋滞	23
3 -(2) Web 構造の変化	24
3 -(2)-1 最近の調査結果から見られた現象	24
3 -(2)-2 原因となる Web 構造の変化の推論	25

図表目次

図表 1	インターネット用途の拡大	5
図表 2	サーチロボット「L o k i」の仕組み	6
図表 3	既知URL 発見率の変化	7
図表 4	クローリング(走査)実施実績	9
図表 5	2004年2月(第11回)調査結果(総サーバ数、総ファイル数、総ページ数、総データ量)	10
図表 6	総ファイル数のファイル種類別シェア(2004年2月調査結果)	11
図表 7	総データ量のファイル種類別シェア(2004年2月調査結果)	12
図表 8	JPドメイン総サーバ数の推移	13
図表 9	主要なセカンドレベルドメイン	13
図表 10	JPドメイン総サーバ数のセカンドレベルドメイン別シェアの推移	14
図表 11	JPドメイン総ページ数推移	15
図表 12	JPドメイン総ファイル数の推移	16
図表 13	JPドメイン総ファイル数の推移(ファイル種類別)	16
図表 14	JPドメイン総データ量の推移	18
図表 15	JPドメイン総データ量の推移(ファイル種類別)	18
図表 16	1万ページあたりの平均データ量と1サーバあたりの平均ページ	20
図表 17	コンテンツ量の推移	21
図表 18	1998年2月を100とする指数で見たコンテンツ量の推移	21
図表 19	汎用ドメインにおける日本語サイトの総コンテンツ量(2002年11月調査)	22
図表 20	L o k i 既知URL 発見率(RK)グラフの2つの現象	24
図表 21	中間的リンク規模のWeb ページの増加	25

はじめに - Webの役割とWebコンテンツ情報量調査の意義

【Webの役割】

1990年代以降、インターネット利用は世界的に急拡大が続いている。その最大の要因はWorld Wide Web¹（以下「Web」と言う。）利用者とその用途の急拡大にある。インターネットの恩恵を「情報流通の時間・距離・費用の克服」と捉えるなら、その用途の拡大にWebほど貢献したツールはないと言える。

Webは文書だけでなく、画像や音声、動画等を含めたマルチメディア情報を提供することができ、ブロードバンドインターネット需要拡大の礎を築いた。Webは膨大なコンテンツを手軽に利用できる形で蓄積するという機能を有するだけでなく、蓄積された情報がオンラインで利用されることによって、ビジネス、生活、社会活動のあり方そのものの変革も促している点で、その経済・社会に与える影響は極めて大きい。

最近のインターネットにおいては、P2P²やマルチキャスト³などの新たなツールが生まれてきているが、利用者・用途の広範さと言う観点で見れば、Webは将来も相当な期間インターネットの中核を占めるツールであり続けるものと思われる。

【Webコンテンツ情報量調査の意義】

こうしたWebの及ぼす経済・社会的影響に鑑み、勃興期からWebの発展の推移を概観的に記録しておくことは、インターネットの社会的役割の変化、Webが原因となる社会活動の変化を裏付ける上で大変重要な作業である。本調査は、1998年から毎年定期的に国内のJPドメインのWebコンテンツ情報量を推計するもので、今年で7年目を迎える。長期にわたり、定期的に、統計的手法を用いてWebのコンテンツ情報量を測定したという点で、世界的にも例を見ない唯一の調査である。

本稿では、Webの勃興期からブロードバンド化が進んだ今日までのWebコンテンツの発展について、2003年度（第11回）調査の最新調査結果とともにその推移を概観し、最近の調査結果から見られたWeb上の新たな現象について考察する。

¹ インターネットやイントラネットで標準的に用いられるハイパーテキスト・システム。ハイパーテキストとは、文章内にあるテキスト文字列が、さらに別のテキストやファイルにリンクしている文書システムのこと。欧州核物理学研究所(CERN)のTim Berners-Lee氏が所内の論文閲覧システムとして1989年に考案したものを基礎としている。広く一般に公開されたのは1991年のこと。インターネット標準のハイパーテキスト・システムとして1990年代中頃から爆発的に普及し、現在では世界規模での巨大なWWW網が築かれている。インターネットで最も多く利用されるアプリケーションである。

² サーバを介さずにコンピューター同士を直接接続してデータをやり取りする通信形態。

³ 単一のバケットを同時に特定の複数のコンピューターに送信する技術。

図表 1 インターネット用途の拡大

	1970年代	1980年代	1990年代
主たる役割	米政府のツール コンピュータ間の データ交換手段	学者・研究者などの コミュニケーション ツール	経済活動・社会活動・生活 の重要なツール 社会インフラ化
中核のツール	TCP/IP・FTP	電子メール	Web
新たな ネット利用者	政 府	大学・研究関係者など	一般団体・個人
新たな用途	データ通信	メール ニュース	マルチメディア情報検索 (仕事・学業・趣味・商品) 商取引・掲示板・報道・放送 教育・公的手続きなど

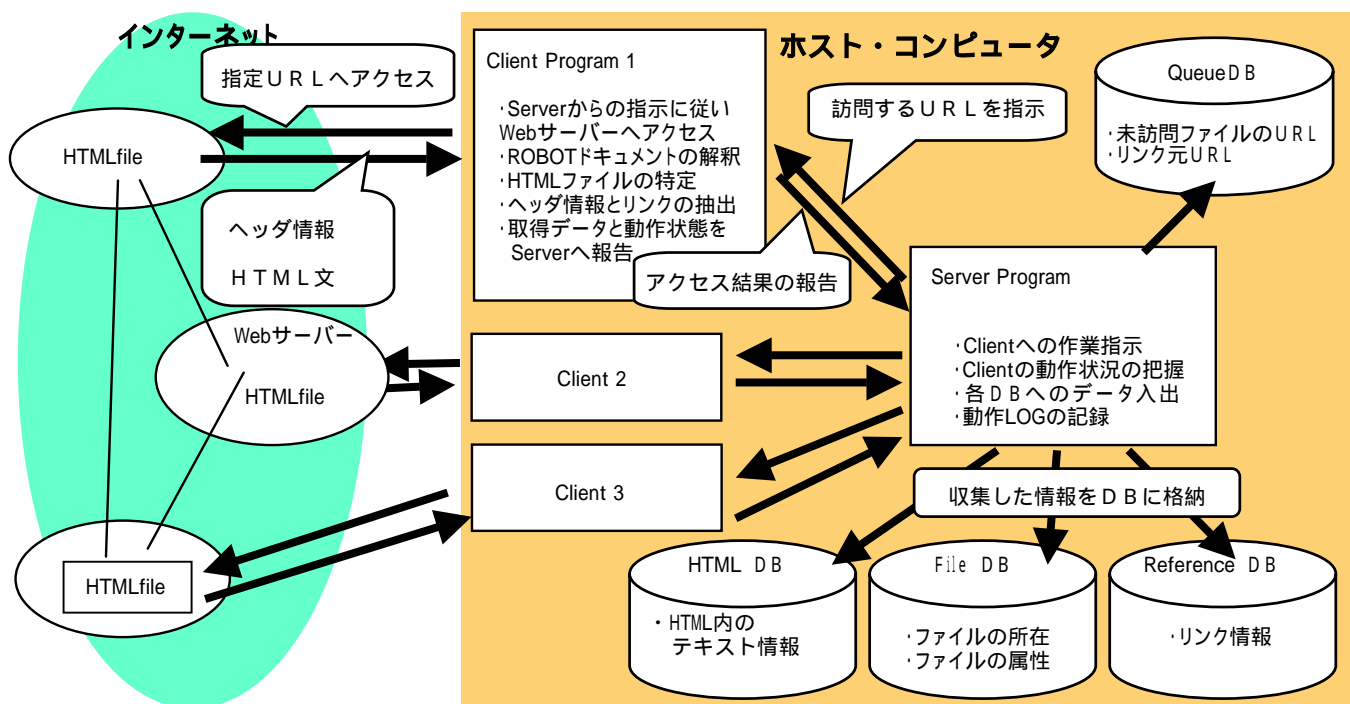
1. 調査の概要

情報通信政策研究所では、旧郵政研究所調査研究部時代の1998年2月から、アライド・ブレインズ株式会社と共同開発したサーチロボットの走査結果をもとに推計する独自の手法により、国内Webコンテンツ量の統計調査を実施している⁴。

Webの膨大な規模や成長速度を考えると、どれほど高性能なサーチロボットを用意しても、公開されているWebページすべてにアクセスすることは現実には不可能である。そこで、Web上に存在する「ある程度の」ページをサーチロボットで調査し、そこで得られるURL(Uniform Resource Location)⁵等のデータをもとに統計的推測を加えて、Webのサーバ数、ファイル⁶数、ページ⁷数、及びデータ量のそれぞれの総数⁸を求めている。

1-(1) 調査の仕組み・原理

図表 2 サーチロボット「Loki」の仕組み



⁴ 「インターネットコンテンツ統計に関する調査研究 郵政研究所月報 2002年9月」

(http://www.soumu.go.jp/iicp/seika/data/research/monthly/2002/168-h14_09/168-asearch2.pdf)

「メディアとしてのWebの成長を測る～サーチロボットを使ったWebコンテンツ統計調査の試み」
マス・コミュニケーション研究 No.62 2003年3月

(<http://www.a-brain.com/HP/rep/rep15/index.html>) を参照。

⁵ インターネット上に存在するHTML(HyperText Markup Language)、文書/データ、画像などのファイルの場所を指し示す記述方式。情報の種類やサーバ名、ポート番号、フォルダ名、ファイル名などで構成されている。

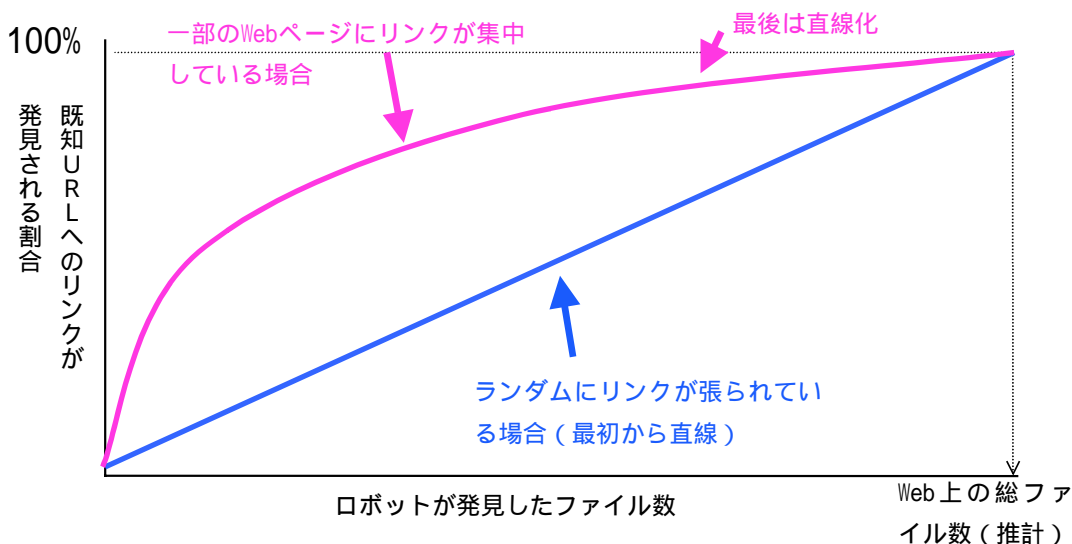
⁶ 本調査では、独立したURLをもつWebサイト上のコンテンツであるHTML、文書/データ、画像、動画、音声などのファイルをいう。

⁷ 本調査では、HTMLファイルと同義とする。HTMLは文書の論理構造や見栄えなどを記述するために使用され、文書の中に画像や音声、動画、他の文書へのリンク情報などを埋め込むこともできることから、HTMLファイルを「ページ」と呼び、他のファイルと区別している。

⁸ 総数の他、ファイル種類別のデータ量、ファイル数を推計している。

本調査のために開発したサーチロボット「L o k i」は、HTMLのハイパーリンクを自動的に辿ってクロール（走査）することにより、インターネット上の広範なWebページのURL情報を自動収集するプログラムである。情報収集のスピードアップのため複数の情報収集プログラム（クライアント）を同時に動かして並列でクロールを行う。サーバープログラムは、これらのクライアントプログラムの動作管理や収集した情報の整理を行う。

図表 3 既知URL発見率の変化



サーチロボットは新しいWebページにアクセスすると、そのHTML文に記述されているリンク情報を取得するが、それらの中にはサーチロボットにとって未知のURLへのリンクもあれば過去既に発見済みのURLへのリンクもある。サーチロボットの走査が進むにつれて、サーチロボットが既知しているURLの数（以下「既知URL」という。）は増えていくので、新たに発見したURLの比率はどんどん低下し、一方で既知URLが重複してみつかると比率はどんどん上昇する。サーチロボットがアクセス可能なすべてのWebページのURLを取得し終わったときには、その中のどのWebページにアクセスしても、そこに記述されているリンク先はすべて既知URLとなるはずである。したがって、サーチロボットがある程度数のURLを取得した段階で、新たに取得したリンク情報の中の既知URLの比率を調べれば、それまでに取得したURLがWeb上のファイル全体の何パーセントにあたるか推定できることになり、すべてのWebページにアクセスしなくてもWebの総ファイル数を推計できることになる。

この手法でWebの総ファイル数等を推計する場合、Webページ間のリンクがWeb全体にランダムに、かつ偏りなく張られていれば、サーチロボットが発見したファイル数と既知URLが発見される割合は比例し、サーチロボットの調査が進むにつれて、既知URL発見率は図表6のように直線的に上昇するはずである。しかし、実際のWebではリンクの分布に大きな偏りがあり、一部の「有名」Webページに多数のリンクが集中している。このような有名Webページは早期にサーチロボットに発見される確率が高く、またいったん発見されると既知URL発見率の値を一気に高めることになる。つまり、実際の調査では、既知URL発見率の値は直線的に上昇するの

ではなく、初期の段階に急激に上昇し、その後漸増し、最後にリンク構造が均一化すると直線的に漸増すると考えられる。このような前提に立ち、グラフが直線になった時点で、線形近似法⁹という手法を用いて Web の総ファイル数の推計を行っている。

1-(2) 調査対象

本調査で対象としている Web の範囲は、一般に公開されている J P ドメイン¹⁰の Web ページである。最近では、国内企業等が公開している「.com」などの汎用ドメイン¹¹をもつものが増えているが、本調査では J P ドメインのみを対象としている¹²。

なお、サーチロボットは必ずしもすべての Web ページにアクセスできるわけではなく、次のような Web ページやファイルは収集ができないため、調査対象から外れている。

- ・ 外部からのリンクを持たないもの
- ・ 会員限定等のアクセス制限のある Web ページ
- ・ 「robots.txt」でサーチロボットのアクセスを禁じている Web ページ
- ・ CGI¹³等のスクリプトで自動生成されるもの
- ・ 検索エンジンで自動生成されるもの
- ・ 再生、停止ボタンが埋め込まれている動画ファイル
- ・ ストリーミングの動画、音声ファイル¹⁴

⁹ 情報通信政策研究所（旧郵政研究所）とアライド・ブレインズ株式会社の共同で特許出願中。

¹⁰ 国ごとに割り当てられる国別ドメイン名（ccTLD）で、日本の ccTLD である「.jp」のこと。セカンドレベルが組織種別をあらわす属性型 J P ドメイン名やセカンドレベルに取得者の希望する名前を登録する汎用 J P ドメイン名などがある。

¹¹ 国の区別なく世界中で自由に取得可能な分野別ドメイン名（gTLD）、「.com」_、「.net」_、「.org」等

¹² J P ドメインを対象としている理由については、後述 3-(1)-1-1 を参照。なお、2002 年 11 月調査は単年で汎用ドメインも調査対象に含めている。

¹³ Common Gateway Interface の略。Web サーバが、Web ブラウザからの要求に応じて、プログラムを起動するための仕組み。従来、Web サーバは蓄積してある文書をただ送出するだけであったが、CGI を使うことによって、プログラムの処理結果に基づいて動的に文書を生成し、送出することができるようにするもの。

¹⁴ Windows Media Player や Real Player のファイルがこれに該当する

2 . 調査結果について

2 -(1) これまでの調査実施状況

サーチロボットを使った調査は 1998 年 2 月の初回調査から 2001 年 8 月の第 8 回調査までは半年ごとに実施していたが (2 月と 8 月)、2002 年以降は年 1 回のペースで実施している。

図表 4 クローリング(走査)実施実績

	サーチロボット走査期間
第1回	1998年 2月10日～ 2月26日
第2回	1998年 8月 3日～ 9月 7日
第3回	1999年 2月16日～ 3月11日
第4回	1999年 8月 4日～ 9月26日
第5回	2000年 1月17日～ 3月 7日
第6回	2000年 8月30日～ 9月27日
第7回	2001年 2月10日～ 3月19日
第8回	2001年 7月20日～10月30日
第9回	2002年 2月 3日～ 4月23日
第10回 ¹⁵	2002年 10月22日～11月 6日
第11回	2004年 1月 5日～ 2月26日

¹⁵ 第 10 回調査ではサーチロボットに株式会社エヌ・ティ・ティ・エックス (現エヌ・ティ・ティレゾナント株式会社) のサーチエンジンを使用。

2-(2) 調査結果

2-(2)-1) 2004年2月(第11回)調査結果

2-(2)-1)-1. JPドメイン総コンテンツ量

【総サーバ数・総データ量の伸び大きく、総ファイル数の伸び低調】

総サーバ数、総ファイル数、総ページ数、総データ量とも前回(2002年11月)調査に比べ増加しており、今回の調査では、総データ量は13,609ギガバイト(GB)¹⁶と推計された。

伸び率で見ると、前回調査との比較では、総ファイル数の伸び(6.39%)に比べ、総サーバ数の伸び(37.99%)と総データ量の伸び(34.08%)が大きい。ファイル数の少ないファイル構造の簡単なサーバが多数出現していると考えられる。また、総ファイル数の伸びに比べて総データ量の伸びが大きいことから1ファイルあたりのデータ量は増えているものと考えられる。

図表 5 2004年2月(第11回)調査結果(総サーバ数、総ファイル数、総ページ数、総データ量)

	2002年 11月(参考)	2004年 2月
総サーバ数(台)	308,000	425,000
前回比伸び率(%)		37.99
総ページ数(万P)	7,438	8,590
前回比伸び率(%)		15.49
総ファイル数(万F)	27,421	29,173
前回比伸び率(%)		6.39
総データ量(GB)	10,150	13,609
前回比伸び率(%)		34.08

¹⁶ バイトは情報量を表わす単位で1ビット(bit)を8つ集めた情報量。1ビットで0か1かの2種類の情報を表わす。1,024バイトを1KB(キロバイト)、1,024KBを1MB(メガバイト)、1,024MBを1GB(ギガバイト)、1,024GBを1TB(テラバイト)と呼ぶ。

2-(2)-1)-2. JPドメイン総ファイル数のファイル種類別シェア

【ファイル数では画像が多く、次いでHTML】

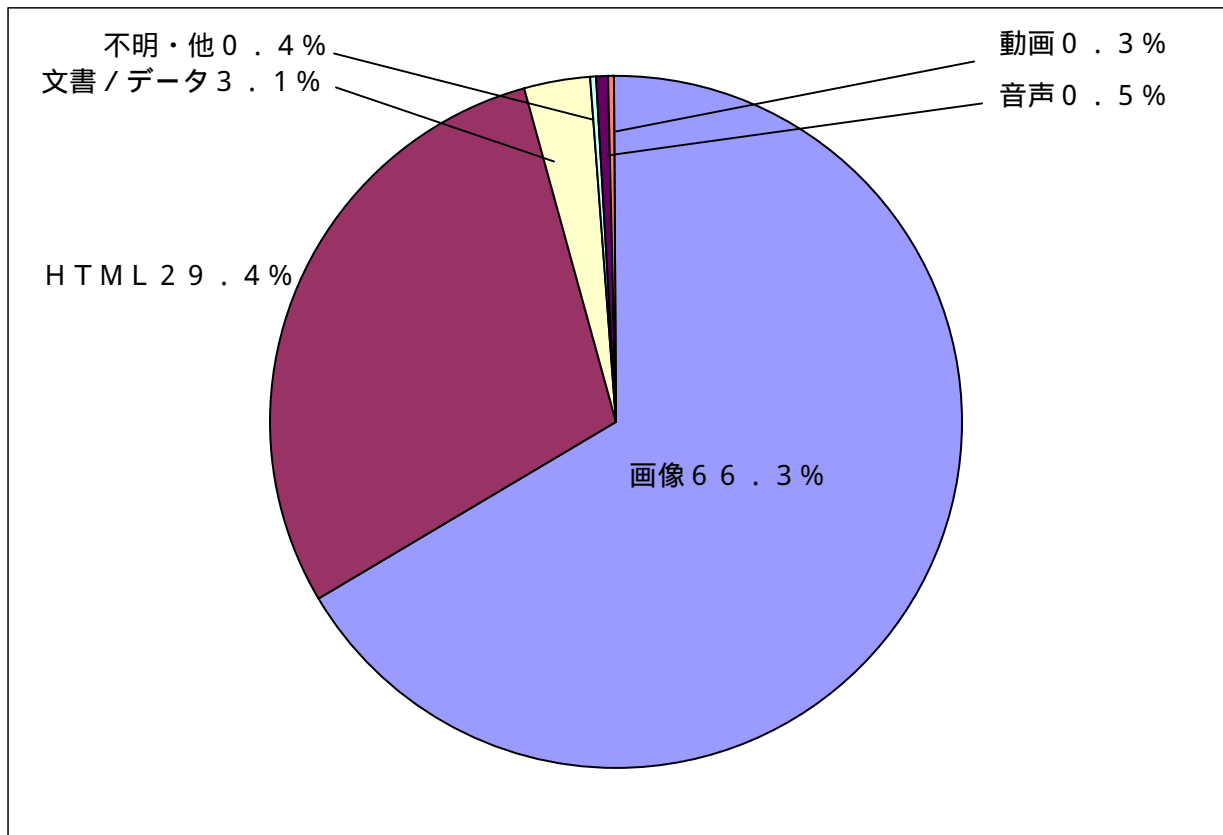
総ファイル数のファイル種類¹⁷ごとの構成比で見れば、多くは画像であり、全体の66.3%を占める。次いでHTMLが29.4%となっている。

動画、音声はそれぞれ0.3%、0.5%を占めるにすぎない。

図表 6 総ファイル数のファイル種類別シェア (2004年2月調査結果)

(単位: %)

画像	HTML	文書/データ	音声	動画	不明/他
66.3	29.4	3.1	0.5	0.3	0.4



¹⁷ ファイルの種類は URL 最後の「拡張子」に基づいて分類している。各種類別の代表的な拡張子は次のようなものである。

- HTML: 「.htm」「.html」
- 画像: 「.jpg」「.gif」「.bmp」「.pict」「.tif」「.eps」「.png」
- 動画: 「.mpg」「.avi」「.mov」
- 音声: 「.au」「.ra」「.midi」「.mp3」「.rmi」「.wav」
- 文書/データ: 「.pdf」「.txt」「.doc」「.jwv」「.lzh」「.tar」「.xls」「.exe」「.java」

2-(2)-1)-3 . J P ドメイン総データ量のファイル種別別シェア

【データ量では動画ファイルが3割を占める】

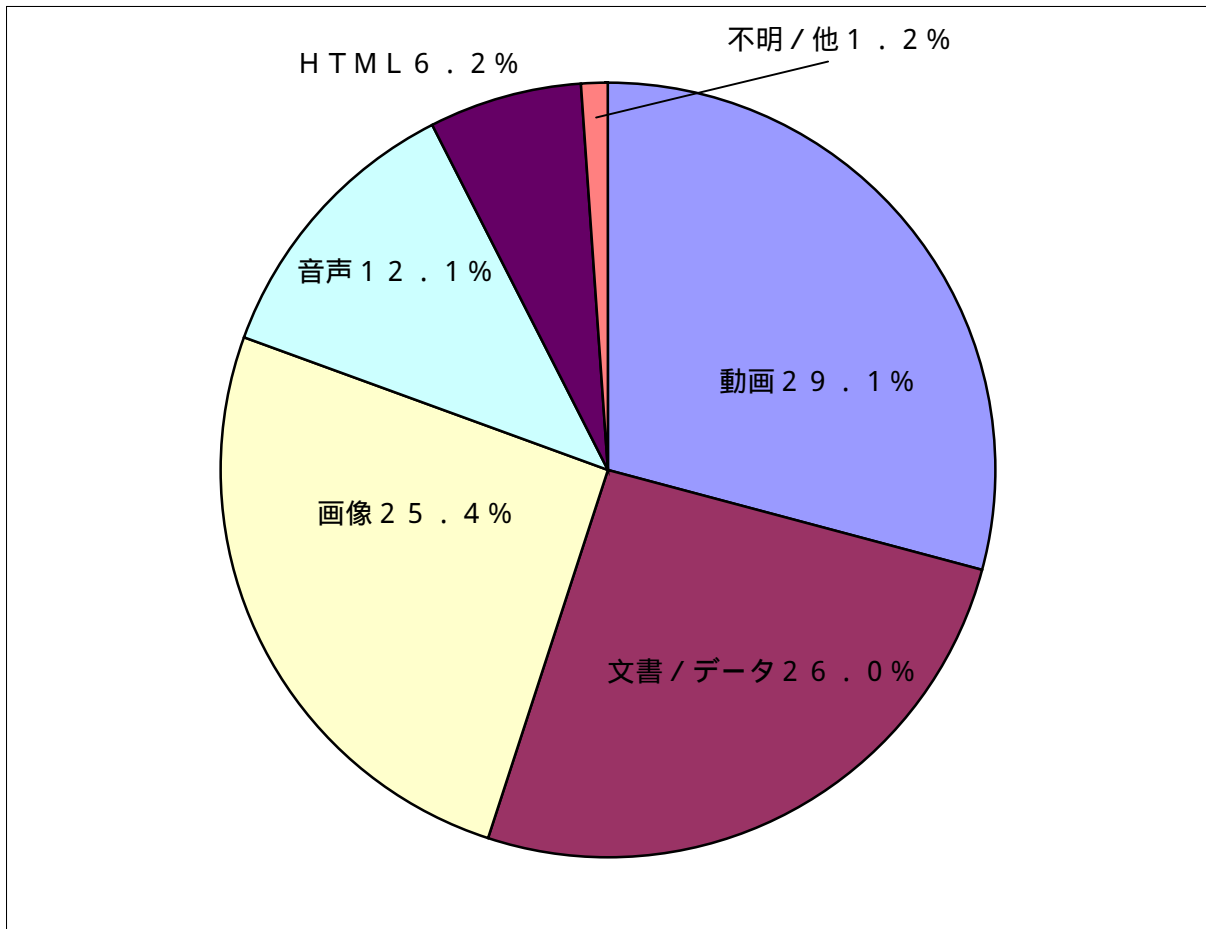
総データ量のファイル種別別構成比で見ると、動画・音声等のマルチメディアファイルや文書・データのデータ量が大きな比率を占め、動画が29.1%と最も多く、次いで文書/データが26.0%、画像が25.4%、音声12.1%となっている。

動画・音声はは総データ量のファイル種別別構成比では約4割を占めるものの、総ファイル数のファイル種別別構成比ではわずか0.8%に過ぎない(図表6参照)。動画1ファイル当たりのデータ量が他のファイルに比べ圧倒的に大きいことが伺える。今後、動画をはじめとするマルチメディアファイルの割合が総ファイル数においても増加すれば、Web上の総データ量はますます増加していくものと考えられる。

図表 7 総データ量のファイル種別別シェア (2004年2月調査結果)

(単位: %)

動画	文書/データ	画像	音声	HTML	不明/他
29.1	26.0	25.4	12.1	6.2	1.2



2-(2)-2) 過去の調査結果との比較

2-(2)-2)-1. JPドメイン総サーバ数の推移

【総サーバ数は一貫して順調な伸び】

JPドメインの総サーバ数は1998年初回調査から一貫して増加しており、特に最近2年間は連続して増加ペースが速まっている。

図表 8 JPドメイン総サーバ数の推移

	1998年		1999年		2000年		2001年		2002年	2002年	2004年
	2月	8月	2月	8月	2月	8月	2月	8月	2月	11月	2月
総サーバ数 (台)	36,000	54,000	75,000	85,000	95,000	120,000	152,000	177,000	197,000	308,000	425,000
前回は伸び 率(%)	-	50.00	38.89	13.33	11.76	26.32	26.67	16.45	11.30	56.35	37.99

【構成比ではacドメインは縮小、coドメインが拡大】

JPドメインサーバのセカンドレベルドメイン別の構成比を見ると、調査開始当初において最も構成比が大きかったacドメインは、その後構成比では縮小を続け、代わってcoドメインの数が増加し、この4年間構成比では全体の過半数を占めるようになってきている。調査初期の頃は、学術機関における情報交換のツールであったWebが、徐々に企業活動のツールにシフトしたことが伺える。

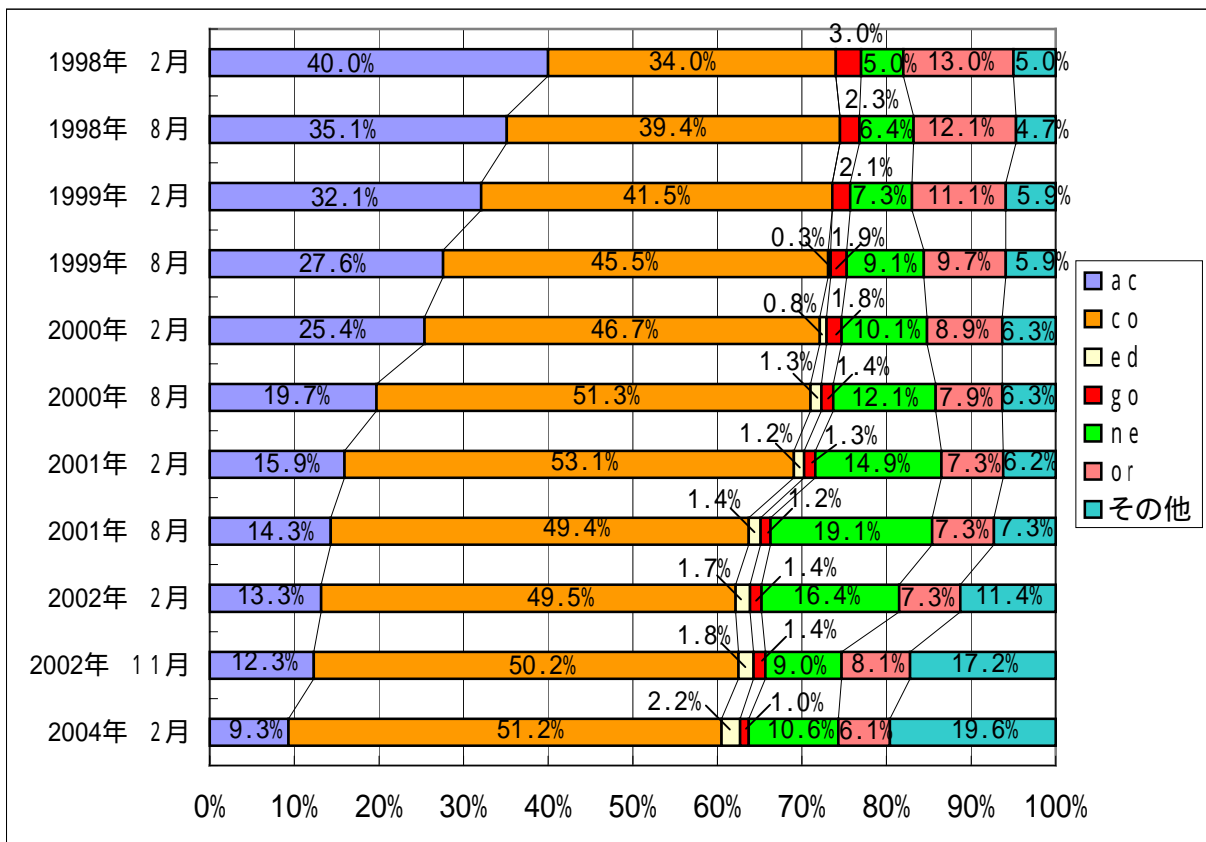
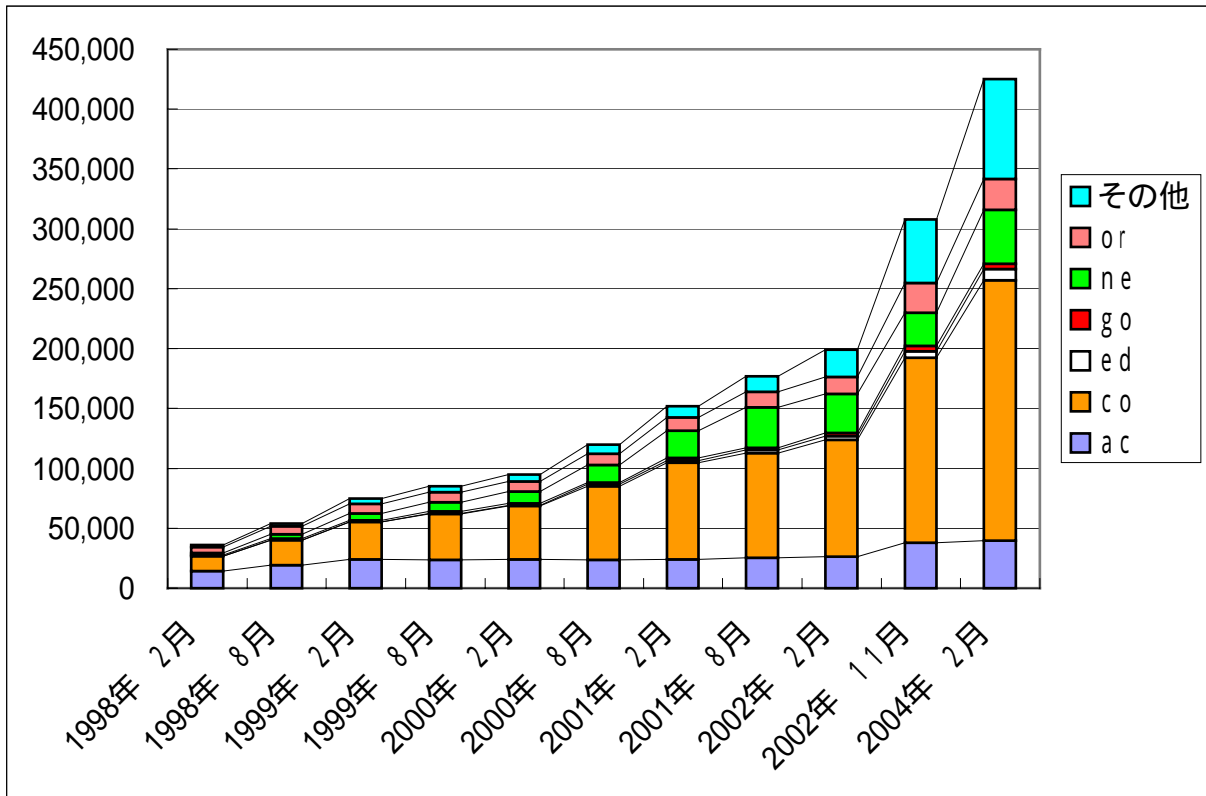
近年の顕著な動きとしては、「その他」ドメインの急速な増加があげられる。「その他」ドメインは2004年2月時点の構成比で約19.6%を占めるところまで来ているが、これは2001年以降の汎用JPドメイン¹⁸の導入・普及の影響が大きいものと思われる。

図表 9 主要なセカンドレベルドメイン

ac・・・大学系教育機関等	go・・・政府機関
co・・・一般企業等	ne・・・ネットワークサービス等
or・・・会社以外の団体等	ed・・・小、中、高等学校

¹⁸ 属性型JPドメイン名は1つの組織につき1つしか使用することができないなどの制約があるが、汎用JPドメインは同一のものがない限りいくつでも自由に登録できる。

図表 10 JPドメイン総サーバ数のセカンドレベルドメイン別シェアの推移



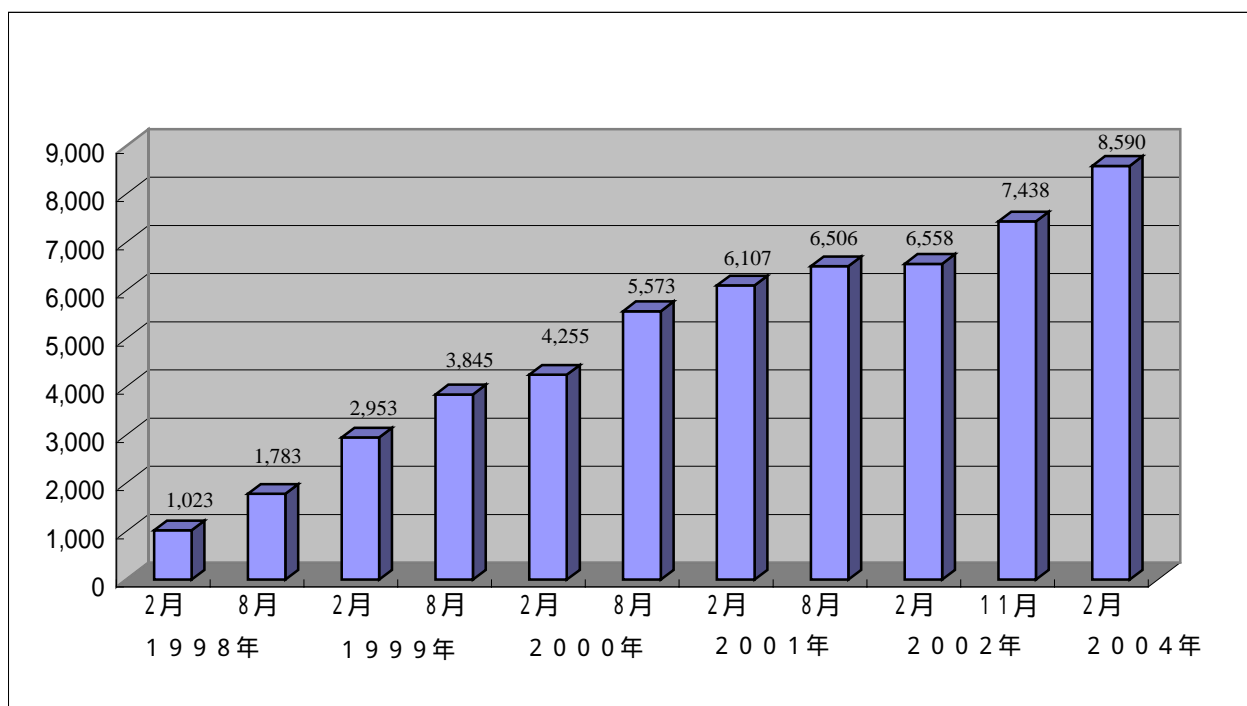
2-(2)-2-2. JPドメイン総ページ数の推移

【総ページ数は一時伸び率が鈍化した近年復調】

JPドメインの総ページ数は、2001年から一時伸び率は低下していたが、この2年間にまた増加ペースを早めている。

図表 11 JPドメイン総ページ数推移

	1998年		1999年		2000年		2001年		2002年	2002年	2004年
	2月	8月	2月	8月	2月	8月	2月	8月	2月	11月	2月
総ページ数 (万P)	1,023	1,783	2,953	3,845	4,255	5,573	6,107	6,506	6,558	7,438	8,590
前回伸び 率(%)	-	74.29	65.62	30.21	10.66	30.98	9.58	6.53	0.80	13.42	15.49



2-(2)-2)-3 JPドメイン総ファイル数の推移

【総ファイル数は2002年後半に急速な伸び】

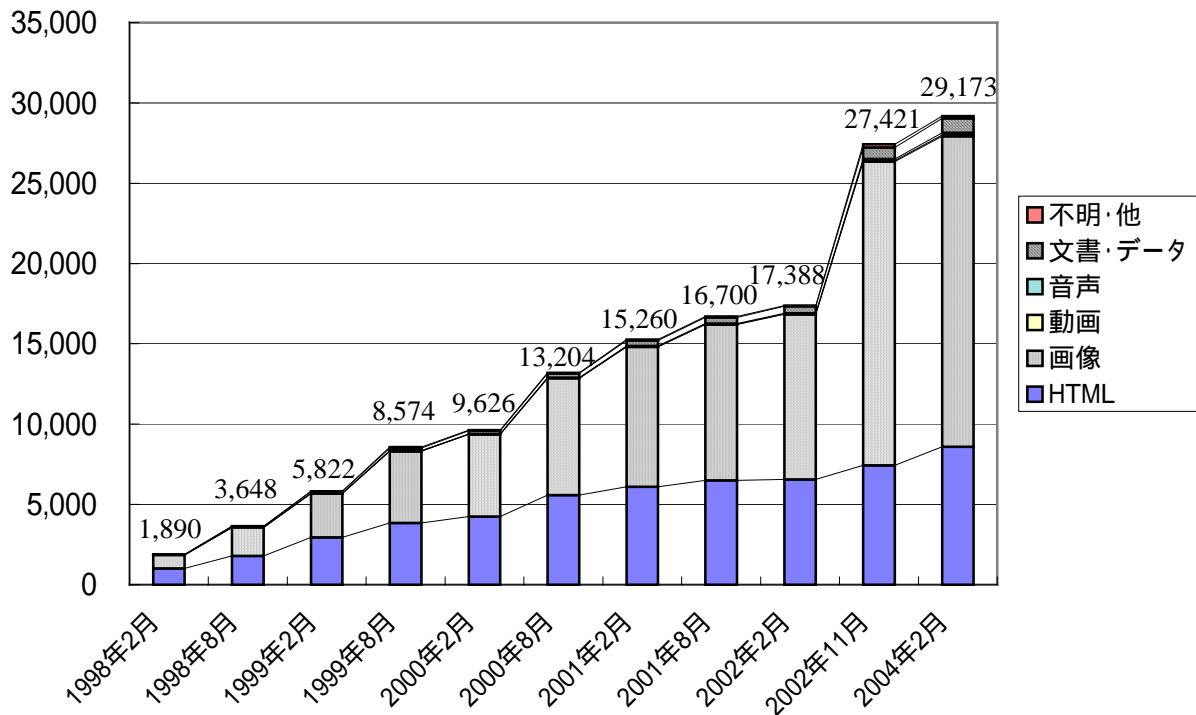
JPドメインの総ファイル数も一貫して増加を続けているが、2001年頃にはその増加率は低下傾向にあった。しかし、2002年後半からは一転して大きな伸びが見られた。

ファイル種別に見ると、画像ファイルが初回調査から高い割合を占めている。

図表 12 JPドメイン総ファイル数の推移

	1998年		1999年		2000年		2001年		2002年	2002年	2004年
	2月	8月	2月	8月	2月	8月	2月	8月	2月	11月	2月
ファイル数 (万P)	1,890	3,648	5,822	8,574	9,626	13,204	15,260	16,700	17,388	27,421	29,173
前回比伸び 率(%)	-	92.91	59.59	47.27	12.27	37.17	15.57	9.44	4.12	57.70	6.39

図表 13 JPドメイン総ファイル数の推移（ファイル種類別）



(単位：万F)

	1998年		1999年		2000年		2001年		2002年	2002年	2004年
	2月	8月	2月	8月	2月	8月	2月	8月	2月	11月	2月
HTML	1,023	1,784	2,953	3,846	4,255	5,573	6,107	6,506	6,558	7,438	8,589
画像	827	1,775	2,727	4,469	5,103	7,277	8,704	9,707	10,288	18,918	19,339
動画	2	4	5	7	8	10	12	12	14	61	81
音声	3	10	11	25	30	34	40	37	36	89	141
文書・データ	25	61	116	173	198	270	348	388	436	720	899
不明・他	11	14	11	54	32	40	49	51	55	195	125
合計	1,890	3,648	5,822	8,574	9,626	13,204	15,260	16,700	17,388	27,421	29,173

2-(2)-2)-4 JPドメイン総データ量の推移

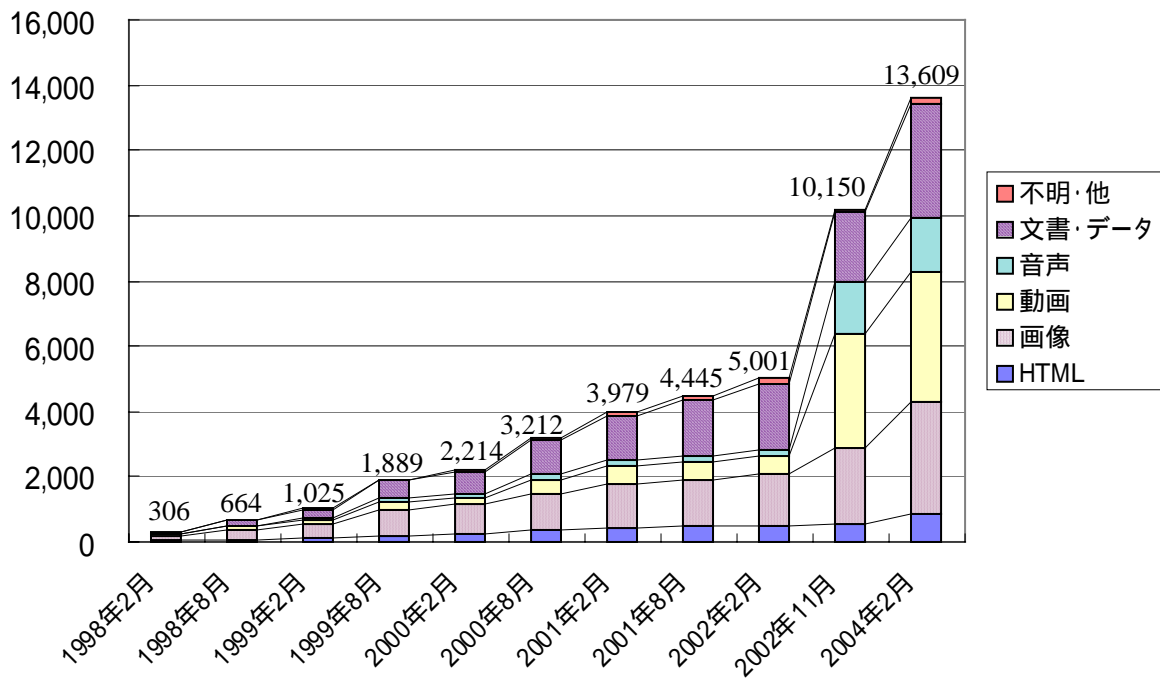
【マルチメディアデータの増大に合わせ順調な伸び】

データ量は、もともとファイルやサーバに比べて伸びが大きい傾向があり、2002年11月調査では前回(2002年2月)調査の2倍という大きな伸びがみられた。今回調査でも前回調査比で34%と高い伸び率となっている。ファイル種別に見ると、動画・音声ファイルのデータの伸びが大きく、これらのマルチメディアデータの増加が全体のデータ量を押し上げている。

図表 14 JPドメイン総データ量の推移

	1998年		1999年		2000年		2001年		2002年	2002年	2004年
	2月	8月	2月	8月	2月	8月	2月	8月	2月	11月	2月
総データ量 (GB)	305	664	1,024	1,889	2,214	3,212	3,980	4,445	5,002	10,150	13,609
前回比伸び率 (%)	-	117.70	54.22	84.47	17.20	45.08	23.91	11.71	12.51	102.92	34.08

図表 15 JPドメイン総データ量の推移 (ファイル種別別)



(単位：GB)

	1998年		1999年		2000年		2001年		2002年	2002年	2004年
	2月	8月	2月	8月	2月	8月	2月	8月	2月	11月	2月
HTML	46	86	150	211	256	354	411	468	498	564	846
画像	141	306	409	745	885	1,135	1,364	1,440	1,579	2,317	3,461
動画	40	78	113	280	206	434	540	530	543	3,507	3,962
音声	11	29	39	88	119	155	210	211	216	1,575	1,642
文書・データ	53	151	300	546	709	1,057	1,356	1,677	2,009	2,174	3,540
不明・他	15	14	14	19	39	77	98	119	156	13	158
合計	306	664	1,025	1,889	2,214	3,212	3,979	4,445	5,001	10,150	13,609

2-(2)-2)-5 JPドメイン1万ページ当たり平均データ量とJPドメイン1サーバ当たり平均ページ数の比較

【サーバ1台当たりの平均ページ数は減少、ページ当たりのデータ量は増加】

Webページの閲覧者がダウンロードすることになるデータ量として、1万ページ当たりの平均データ量を比較すると、1998年2月調査では約0.3ギガバイト（GB）であるが、2004年2月調査では1.58ギガバイトと約5.3倍に増加している。現在のWeb閲覧者は6年前より5倍以上のデータをダウンロードしていることになる。

なお、発信者側のサーバに蓄えられる情報量として、サーバ1台当たりの平均ページ数を見ても1998年2月調査の283.3ページから2000年2月調査の447.4ページへと1.58倍まで順調に拡大した後、一転して減少を続けており、2004年2月調査では202.1ページと、1998年2月調査の数字を下回っている。すでに立ち上がっているサーバのページ数が大きく変動するとは考えにくいので、新しく立ち上がるサーバの中にページ数の少ないサーバが多数存在し、それが平均ページ数を押し下げているものと思われる。

しかし、ページにリンクされる平均データ量が5.3倍になっていることを考えると、ページ数としては減少しても、データ量として見ればサーバの情報量は着実に増大していると言える。

図表 16 1万ページ当り平均データ量と1サーバ当たり平均ページ

	総ページ数 (万ページ)	1万ページ当り平均 総データ量(GB)	総サーバ数 (台)	1サーバ当たり平均 ページ数(ページ)
1998年2月	1,020	0.3	36,000	283.3
1999年2月	2,950	0.35	75,000	393.3
2000年2月	4,250	0.52	95,000	447.4
2001年2月	6,101	0.65	152,000	401.4
2002年2月	6,555	0.76	197,000	332.7
2002年11月	7,438	1.37	308,000	241.5
2004年2月	8,590	1.58	425,000	202.1

2-(2)-2)-6 J Pドメイン総コンテンツ量発展推移(年別)

初回の1998年2月調査の推計値を100とし、総サーバ数、総ページ数、総ファイル数、総データ量の各年別増加状況を見たものである。

【総データ量は初回調査の44.6倍】

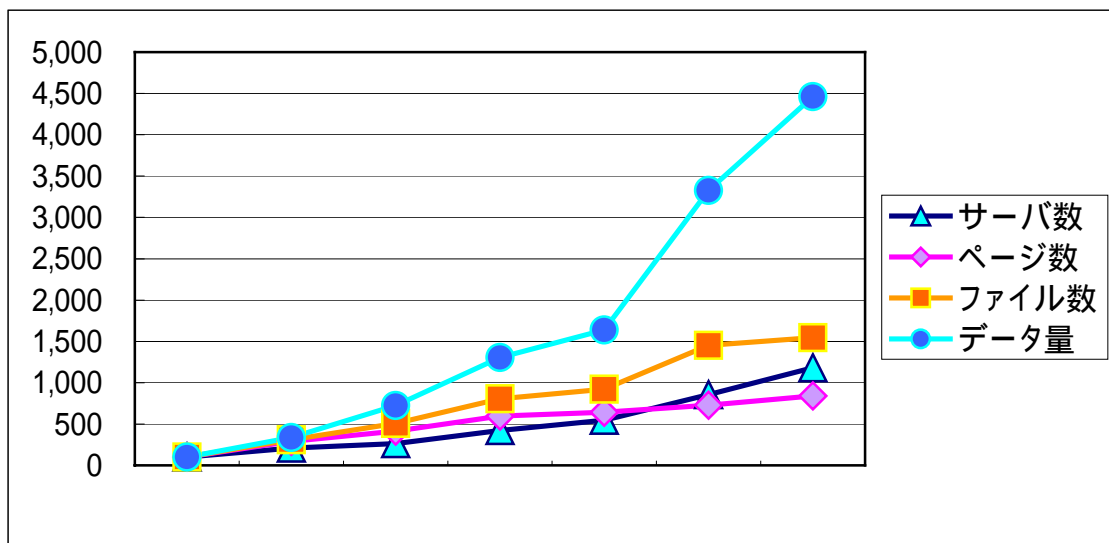
増え方は従来は総サーバ数より総ページ数の方が上回っていたが、近年は総ページ数の伸びが緩やかになり総サーバ数の伸びがあがってきたため、増え方が逆転している。総ファイル数は緩やかな上昇傾向にあるが、総データ量は活発な伸びを示しており、2004年2月は1998年2月に比べ、約44.6倍にまで増加した。

図表 17 総コンテンツ量の推移

推計値	1998年2月	1999年2月	2000年2月	2001年2月	2002年2月	2002年11月	2004年2月
総サーバ数(台)	36,000	75,000	95,000	152,000	197,000	308,000	425,000
総ページ数(万P)	1,023	2,953	4,255	6,107	6,558	7,438	8,590
総ファイル数(万F)	1,891	5,822	9,626	15,260	17,388	27,421	29,173
総データ量(GB)	305	1,024	2,214	3,980	5,002	10,150	13,609

図表 18 1998年2月を100とする指数で見た総コンテンツ量の推移

指数	1998年2月	1999年2月	2000年2月	2001年2月	2002年2月	2002年11月	2004年2月
総サーバ数	100	208	264	422	547	856	1,181
総ページ数	100	289	416	597	641	727	840
総ファイル数	100	308	509	807	920	1,451	1,544
総データ量	100	336	726	1,305	1,640	3,328	4,462



3 . 最近の調査結果から見られた新たな現象

1998 年から国内の Web サーバのコンテンツ量を推計し、定期的に公表している本調査の推計値は、インターネットの規模を表す代表的な指標として広く利用されているところであるが¹⁹、最近の調査結果から Web の規模の拡大や構造の変化が原因と思われる様々な現象が見られた。調査開始から 6 年が経過し、調査開始当初には想定していなかった新たな Web 上の現象について考察する。

3 -(1) 調査設計に関する変化

3 -(1)-1) J P ドメインの国内シェアの変化

本調査は、「勃興期にあった日本国内の Web の発展の軌跡をきちんと記録に留めておくこと」を目的に、国内にある Web サーバのコンテンツ量を推計するものとして始められたため、調査対象を J P ドメインに限定している。

1998 年当時としては、汎用ドメインの国内サイトは特殊で限られたものであったことや、汎用ドメインを有する Web サーバについては、地理的所在を示す情報がとれないので、日本語かどうかは推測できても、国内にあるのかどうか分からない、といった問題があったためである。

しかしながら調査も 6 年を経て国内の Web の状況も確実に変わりつつあり、2002 年 11 月に単年で調査した汎用ドメインを有する日本語サイトの総コンテンツ量（2002 年 11 月）（図表 19 参照）を見ると、汎用ドメインを有する日本語サーバの推定総数は、465,000 台と J P ドメインサーバの推定総数 308,000 台を上回っており、汎用ドメインに海外設置サーバの日本語サイトが多数含まれていることを考慮しても、国内に設置された汎用ドメインサイトも相当数に上ると考えられる。

このような状況の中で我が国の Web コンテンツ量の動きを捉えるには、J P ドメインのみを調査対象としていることについても見直しをする必要があるか、その場合、仮に対象を広げるにしても、汎用ドメインの日本語サイト中、国内のサーバに設置されたものであるかをどのように確認するのか、というようなことについて検討していくことが必要である。

図表 19 汎用ドメインにおける日本語サイトの総コンテンツ量（2002 年 11 月調査）

	サーバ数 (台)	総ページ数 (万 P)	総ファイル数 (万 F)	総データ量 (GB)
汎用ドメイン中日本語によるサイトの 総コンテンツ量	465,000	2,502	9,642	5,444
J P ドメインの総コンテンツ量	308,000	7,438	27,421	10,150

¹⁹ 「IT 戦略会議ベンチマーク集 平成 15 年 12 月 18 日」<http://www.kantei.go.jp/jp/singi/it2/dai22/22siryou6.pdf>

3-(1)-2) ロボットで取得できないマルチメディアデータの増加

サーチロボットはすべての Web ページにアクセスできるようにはなっておらず、次のような Web ページやファイルのデータは収集できない。

- ・ CGI 等のスクリプトで自動生成されるもの
- ・ 再生、停止ボタンが埋め込まれている動画ファイル
- ・ ストリーミングの動画、音声ファイル

こうしたページやファイルは近年 Web 上で主要なツールとして普及してきており、動画や音声についてはむしろこの様な形態のものが主流となっているようである。従って本調査では動画・音声の急激な増加が観測されてはいるものの、情報量推計という観点からは動画・音声の主要な部分が把握できていないという問題がある。

3-(1)-3) Web 規模の増大と調査インフラ増強の問題

Web の規模は年々拡大を続けており、短期間に多くの Web ページにアクセスできるよう Loki も年々改良しているが、そのような対応についても技術的にもコスト的にもおのずから限界がある。

これまで、Loki のプログラムを並列化することにより、複数のマシンを同時に稼働させて各マシンに適切に不可分散を行うなどの改良を行ってきたが、現実には一定数以上の並列処理になると様々な不具合が生じているところであり、この手法による増強にも限界がある。

3-(1)-4) ネットワークのトラフィック増大によるハブの渋滞

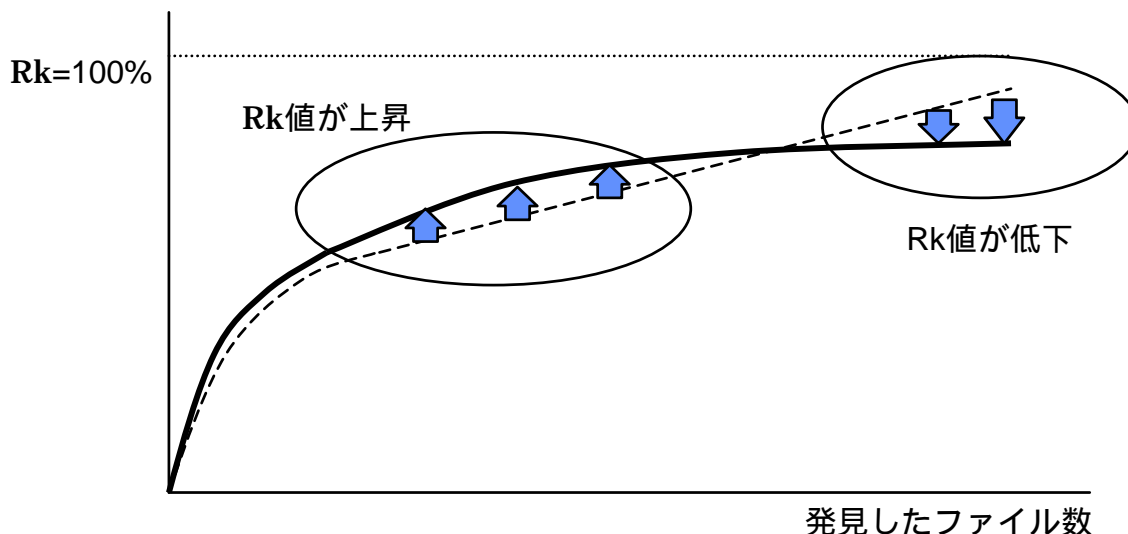
Loki の走査は、ネットワーク環境に左右されるところ、近年調査では、トラフィック増大によるハブの渋滞が走査期間を長引かせる原因となっている。マシンのみならず、回線能力を主とする設置環境の改善（なるべくバックボーンに近いところからインターネットに入るなど）も重要な課題となっている。

3-(2) Web構造の変化

3-(2)-1) 最近の調査結果から見られた現象

1-(1) 調査の仕組み・原理のところでも述べたように、本調査は既知 URL の発見率 (Rk 値) のグラフをもとにファイル総数を推計しているが、そのグラフに関して、最近の調査結果では従来見られなかった 2 つの現象が生じている。

図表 20 Loki 既知URL発見率 (Rk) グラフの 2 つの現象



Rk 値の上昇、曲線化

Rk 値グラフは当初急速に立ち上がり、ある程度までくると徐々に傾きを緩やかにしつつ曲線を描き、最後には傾きが緩やかな直線になってくる。Rk 値グラフの傾きが当初急であることは、取得しているファイルの被リンク数が多いことを、傾きが曲線を描きつつ緩やかになっていくのは、1 日の取得ページの平均被リンク数が絶えず減少しつつづけていることを、最後に傾きが緩やかな直線になるのは、被リンク数が落ちるところまで落ち、疎な状態で均質化していることを示す。

これが調査開始当初に描いていた Web 構造のモデルであり、調査開始当初からこれを裏付ける Rk 値のグラフを得ていたところであるが、最近の調査結果では Rk 値のグラフが急速に上昇し、その後曲線が緩やかになってもなかなか直線化せず、長い期間曲線を描きつつける現象が生じている。今回の調査でもはっきりした直線が得にくく、直線の判断が極めて難しい状況であった。

調査終盤の Rk 値の横這い化

本調査の推計モデルでは、Rk 値のグラフは最後に緩やかな傾きを持った直線となることが重要であり、この傾きに沿って既知リンク率 100% まで延長線を引くことで推定総数が求められる。これまでの結果では、これを裏付ける結果があらわれていたが、直近の結果では、上記 のようにはっきりとした直線が出にくくなった上、直線化した後にさらに横這いを続けるという現象が起きている。これは実績としての取得ファイル数は増えているにもかかわらず、Rk 値の上昇が止

まっていることを意味する。

3-(2)-2) 原因となる Web 構造の変化の推論

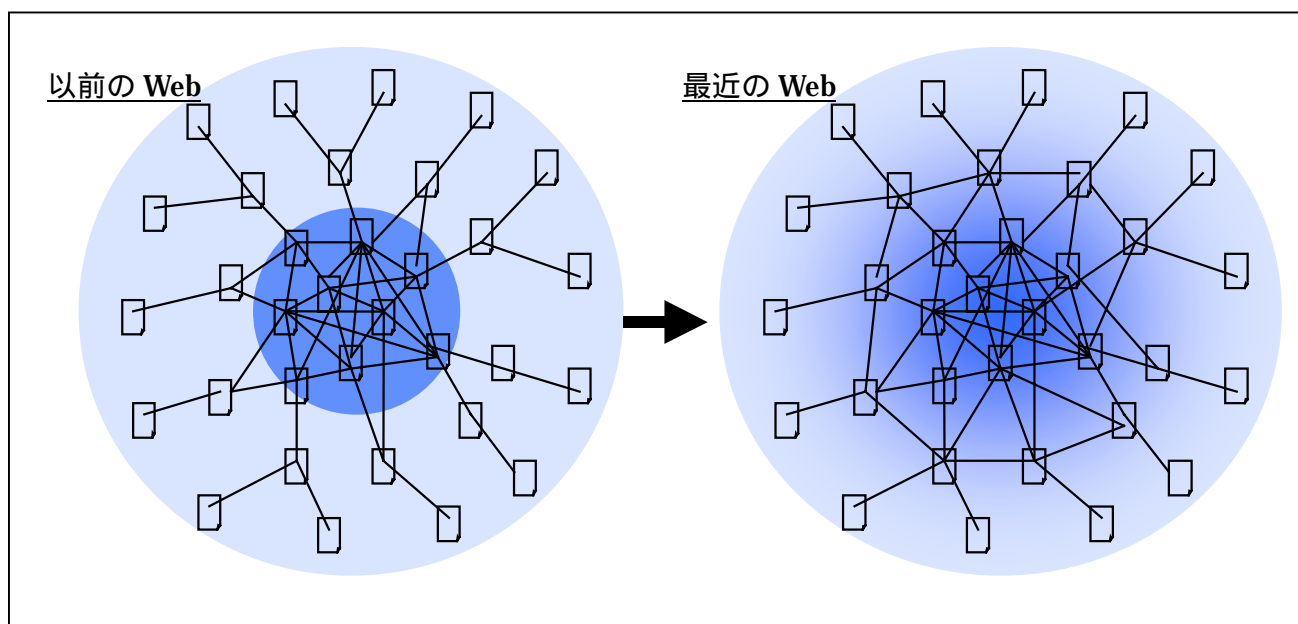
3-(2)-1)で挙げた、最近の調査結果で見られた現象について、その原因を推論する。

【中間的リンク規模の Web ページの増加】

調査開始当初想定していた Web 構造のモデルは、リンクの多く貼られている Web ページと、そうでない Web ページとが 2 極化しているという前提に立ったものであるが、Rk 値が急速に上昇したあと曲線が緩やかになってもなかなか直線化せず、長い期間曲線を描き続ける現象が見られたことから、最近はこのような 2 極化をしておらず、中間的な、様々な被リンク数の Web ページが多数出現して、疎で均質なページとの境界があいまいになってきているものと思われる。

これは、Web の利用が一般にまで広く普及するとともに、その本来的な特徴ともいえる情報を「つなく」という機能が一般の Web ページにも広く浸透し、Web 構造の成熟化の進んでいることを示していると考えられる。

図表 2 1 中間的リンク規模の Web ページの増加



【デッドリンクの増加】

調査終盤の Rk 値の横這い化現象については、Web 周辺部でデッドリンク（リンク先が存在しない、あるいはロボットが辿れないリンク）が増加したためではないかと考えられる。これまでの推計の前提は、デッドリンクは Web 全体でせいぜい数%程度であり、Rk 値の推移に大きな影響は与えないというものであった。しかし周辺部へ行くほどデッドリンクが増加し、かなりの割合でデッドリンクが存在するなら、Rk 値のグラフは 100% に向かわなくなる。

デッドリンクを生む原因としては、長年メンテナンスされていない古いページの増加、アクセス制限をかけているページの増加、自動生成ページの増加が考えられる。

また、携帯電話向けサイトなど、ブラウザ種別に応じてコンテンツを提供するページの増加も原因と考えられる²⁰。

²⁰ Web 空間内にある「勝手サイト」と呼ばれる携帯向けサーバについては、Loki は携帯用のブラウザを有しておらず、勝手サイトからはアクセスを拒否され、デッドリンクとしてカウントしてしまう。