

Web デザイン特別プログラム教育のための小規模サーチエンジンの試作

Trial Development of Small-scale Search Engine to Educate Undergraduate Students in "Web Design" Special Program

幸谷智紀* 竹口友大**

Tomonori KOUYA* Tomohiro TAKEGUCHI**

Abstract: In April of 2009, we will start to educate "Web Design" as a important part of some "Special Programs" in newly started "Faculty of Comprehensive Informatics." In this program, we will reserve some amount of time to spend on learning scripts languages like PHP and Perl and SQL to manipulate databases by project-based learning method. Lack of educational materials for these purposes is currently making us develop such materials. In this paper, we describe brief summary of our small-scale search engine as one of educational materials for "Web design" special program.

1. 初めに

静岡理工科大学総合情報学部では 2009 年度から特別プログラムの一つとして「Web デザイン」を据えることを決定している⁴⁾。いわゆる Project-Based Learning(PBL) の手法を用いて実施されるもので、総合情報学部に属する意欲あふれる学部 2 年生 20 名程度を対象に、World Wide Web(以下 Web と略記)の技術全般を、実践を通じて習得することを目的とするものである。現在の計画では一年間に渡って 168 コマ(1 コマ=90 分)を確保し、5~6 名の教員が Web に関する技術の習得と制作物の評価を行うことを予定しており、我々はこの中で「データベースと Web」、即ち DBMS と Web を組み合わせて利用する 3 層 Web システムを構築する時間を担当することになっている。計画当初はこの時間内にサーチエンジン全体を構築する実習を行う予定であったが、その後の担当教員グループでの話し合いの結果、あくまでその基礎部分を教えるに留めたほうがよいということになった。従ってサーチエンジンの構築そのものは、意欲と能力のある学生向けのオプション教材、および卒業研究のための教材という位置づけである。

教育としてサーチエンジンを構築する試みはテネシー大学で既に行われコースウェアが出版されている³⁾が、これは線型計算と情報技術との橋渡しを目指したものであるのに対し、我々の目指すのは 3 層 Web プログラミング技術の応用としてのサーチエンジンシステムである。従って、基本的には既存のスクリプト言語の習得と、そこから SQL を用いて DBMS を操作するためのプログラミング実習が主となるが、それらを予備知識のない学生に習得させるための教材は現在適切なものが存在しない。そのため、我々は学内 LAN の環境でも設置可能な小規模サーチエンジンシステムを作り、それを段階を追って技術を習得しながら作り上げるための教材²⁾もあわせて執筆した。2007 年度の幸谷研究室における情報セミナー II ではこれを用いて 3 層 Web プログラミングの技術習得を行っている。最終的にはこのシステムを様々な用途に応用していくことを予定しており、従って、用途に合わせた実践的なサーチエンジンシステムとその技術を学ぶための教材開発は今後も継続していかねばならない。そのため、プロトタイプ版でも現状の商用サーチエンジンの基本機能を持つものを作り上げ、その性能評価を行ってシステムの限界を知っておく

必要がある。

本稿ではまずこのサーチエンジンシステムを構築するに至った経緯を説明する。次に現状の商用サーチエンジンシステムの例として Google の機能を解説する。そしてこれを土台にして構築した我々の小規模サーチエンジンシステムの解説を行い、その性能評価を行った結果を示す。最後に、まとめと今後の課題を提示して、次年度以降の開発計画への指針としたい。

2. Web デザイン特別プログラムと 3 層 Web プログラミング教育

Web デザイン特別プログラムは総合情報学部に属する 5 名の教員(菅沼・宮岡・金久保・手島・幸谷)に非常勤講師を加えた 6 名で担当し、主として人間情報デザイン学科の Web デザインコースに進む 2 年生が受講することを想定している。

このプログラムの目的は「Web デザイン」を総合的に学ぶことである。近年の Web は静的な HTML ファイルによる Web ページだけでは語れず、動画ファイルや Flash などを用いて視覚効果を高めたり、バックエンドにデータベースを用いる 3 層 Web システムを用いることが当たり前になっている。当然、TCP/IP Protocol Suite を土台とした Web サーバシステムの運用についてもそれなりの知識が必要となる。このような現代の Web を成立させるための IT 技術を総合的に学び、加えて狭義のデザイン能力とそこから受ける心理的効果についても学び、身に付けることが本プログラムの目的である。

このように様々な異分野の技術や知識を学ぶため、一年間に渡り、毎週 6 コマを使い、Project Based Learning の手法を用いて、理論と実践を交互に組み合わせながら Fig.1 に示すような順序で総合的な Web デザイン能力を磨いていく。

我々はこのうち「データベースと Web」36 コマ分を担当する。この部分のシラバス案の詳細を Fig.2 に示す。ご覧の通り、この中では 3 層 Web プログラミング技術を中心に学んでいくことになる。

ユーザが直接操作する Web ブラウザ、Web サーバ、バックエンドの DataBase Management System(DBMS)の 3 層から成る Web システムを構築するためのプログラミングを、本稿では「3 層 Web プログラミング」と呼ぶことにする。この歴史は古く、CGI は Web 草創期には利用されていたので、それと同じ時期には考案されていたのではないかと推察される。既に 1996 年にはまとまった邦書⁵⁾が出版されているが、10 年以上経った現在でも、フリーな DBMS が主流になってきたと

2008 年 3 月 4 日 受理

*理工学部情報システム学科

** (資) わいにじ

特別プログラム概要
HTML
JavaScript
PhotoShop & Flash
動画コンテンツの作成
Webページの心理評価(1)
Webデザインとユーザビリティ
Semantic Web/インターネットビジネス/情報倫理
Webページの心理評価(2)
前期報告レポート作成
2年生前期

9月	3DCGとWeb3D
10月	Web3Dの制作
11月	データベースとWeb
12月	Java
	Webデザイン
1月	Webページの心理評価(3)
	最終発表会と全体講評
	2年生後期

Fig. 1: Web デザイン特別プログラムのシラバス案

いう以外の技術の変化はない。それだけ「枯れた」技術と言えよう。この枯れた技術の現代的応用としてはサーチエンジンというものが面白いのではないかと我々は考えている。

しかし、現状では学生のプログラミング技術の習得状況は甚だ怪しく、プログラミング以外の関連事項も含め、一通りのことを PBL の中で教え込む必要があると思われる。以下、3 層 Web プログラミング教育に必要な重点事項を述べる。

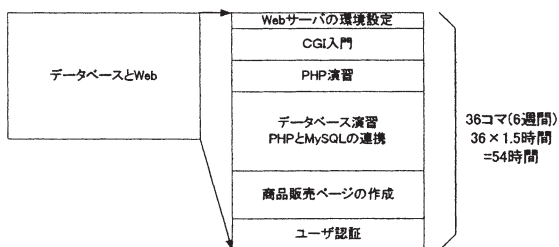


Fig. 2: 「データベースと Web」のシラバス案

2.1 Web サーバの運用技術

インターネットにさらされる環境で運用することを前提とした、Web サーバの運用技術の習得は必須と考えている。最近では ISP が用意するお仕着せの環境で CGI を動かすことが当たり前のことになっているが、じっくり時間をかけて学ぶことのできる機会に、システムの根本に触れておくことは無駄ではない。また、仮想マシン環境が一般化してきている昨今の現状を考えると、管理者 (root) としてサーバに触れる機会は逆に増えてくるとも考えられる。

そこで、昨年整備した VMware 環境⁶⁾ を利用し、OS(CentOS) のセットアップから Web サーバの環境設定まで全てを学生個人で行ってもらい、その環境の上で全ての実習を行うようにしようと考えている。Web サーバは OS に標準で搭載されている Apache を使用し、httpd daemon の起動と停止、httpd.conf の設定変更、アクセスログの読み方とログ解析ツールの利用法等を習得させる。UNIX の CUI に不慣れな学生が殆どだが、我々の経験上、GNOME のような GUI も併用すれば、実習に

支障をきたすことはないと思われる。

2.2 Perl/PHP の知識

3 層 Web プログラミングは、スクリプト言語だけではなく、Java や C/C++ といったコンパイラ言語でも可能である。最近ではむしろ後者、特に Java による実装が増えているように思われる。Multi-core CPU が増えつつある昨今では、Multi-thread 化の容易な言語で実装した方がパフォーマンスの点で有利になることが多いだろう。

しかし、Perl による CGI の実装は長い間標準として使用されてきており、現在でも Movable Type⁸⁾ のような主要な Web アプリケーションでも Perl を利用している例は少なくない。また、後述するように、我々の実装した Web ロボットは、長年の実績を積んできた Perl Module(libwww, Encode 等) を多数利用し、行数を抑えている。将来的には Java 等のコンパイラ言語への移行も考える必要があるが、現状、我々が持っているネットワーク資源を考えると、Perl で十分という面がある。

そこで我々は、Perl による CGI やデータベースプログラミングを教えると共に、並行して PHP についても触れる予定でいる。言語に関わらず、データベースの操作は SQL を発行してデータを操作するだけ、という基本を知る意味でも複数の言語に触れる機会があった方が良い。また、Perl に似ていながら可読性の良いスクリプトが書け、HTML に埋め込む形で容易に変更が可能な言語が存在することも知っておいた方が良い。

2.3 データベース (SQL) の基礎知識

コンピュータを専門とする大学学部において、データベースを講義しないところは皆無であろう。本学でも開講されているが、受講した学生の様子を見ると、ほとんど身につけていないというのが現状である。これは担当教員の責任というよりは、座学として RDB 理論を教えるという教育体制自体に無理があると思われる。実際、履修した学生に聞いてみると、RDB の操作は集合とその演算として考えられるという基本は理解しているようだが、それを SQL 文として書き下して活用することには繋がっていないようである。

従って、じっくり実習期間が確保できる PBLこそ、データベース技術の習得にふさわしい場と言える。前述したように、

データベースと Web の連携は枯れた技術であるが、現在のよう
にフリーで使用できる DBMS が存在し、外部ストレージが
Tera バイト単位で扱えるようになった昨今でこそ、その真価
を個人で体験できるようになったと言える。

特に今回実装するサーチエンジンでは、扱うデータ量は膨
大になる。また、昨日までアクセスできていた Web サイトが
今日は存在しない、ということも日常茶飯事である。収集し
たデータの更新・削除を高速におこなうためには、信頼性の
高いトランザクション処理が可能な DBMS が欠かせない。

従って、我々の教育では DBMS を扱うためのごく基本的な
SQL 文の使い方だけを扱うようにし、データの更新・削除・
検索を行う 1 行 SQL 文が書き下せる水準に達すればよしと
する。

3. Google の検索機能

以上述べてきたように、我々の開発するサーチエンジンは 3
層 Web プログラミング技法を学ぶための教材として適切な規
模のものである必要がある。また、ネットワークやコンピュ
タリソースの限られた環境で動作させるため、集められる URI
数は自ずと限られてくる。従って、サーチエンジンとしては
行数も URI 数も必然的に小規模なものとなる。

とはいえ、小規模ながらも Web デザイン特別プログラム教
育の一環として使用される以上、現代的なサーチエンジンの
基本機能は備えていなければならない。そこで、世界的に最
も成功している Google の機能の主要な機能を備えたものを
目指すことにする。

Google の成功は、特に PageRank と呼ばれる被リンク数に
基づく URI のランク付け機能によってもたらされたと言われ
ている¹⁰⁾。しかし、最近では PageRank に加えて、検索キ
ャワードや URI 間の関連性を柔軟に取り入れた機構も導入さ
れているように見受けられる。ここでは PageRank による URI
のランク付けの機構と、表面的に観察できるそれ以外の機構
について述べる。

3.1 PageRank による重みづけ

PageRank に関しては既に Web 上¹⁰⁾ や優れた書籍¹²⁾ によ
る解説が溢れているので、ごく簡単な解説に留める。基本的
には Fig.3 に示すようにして、被リンク関係に基づいた確率遷
移行列 H を導出することになる。

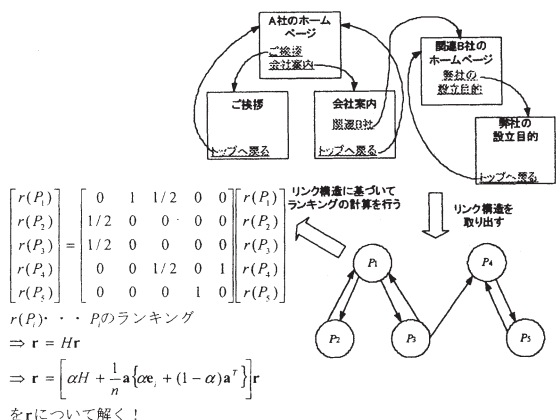


Fig. 3: PageRank 計算のための確率遷移行列の導出

この際、Tree 構造の末端、ループ構造もあり得るので、適度

なランダム移動成分を加味して次のような処理を行って Google
行列 G を導出する。

1. H の列成分が全てゼロの列を $1/n$ で置き換えた S を生成する。もし i 列目が全てゼロであれば:

$$S = H + \frac{1}{n} \mathbf{e} \mathbf{e}^T \quad (1)$$

となる。ここで \mathbf{e}_i は単位ベクトルであり、 $\mathbf{e} = \sum_{i=1}^n \mathbf{e}_i$ である。

2. $0 < \alpha < 1$ という定数 α を決めておき、Google 行列 G を

$$G = \alpha S + \frac{1-\alpha}{n} \mathbf{e} \mathbf{e}^T \quad (2)$$

$$= \alpha H + \frac{1-\alpha}{n} \mathbf{e} (\alpha \mathbf{e}_i^T + (1-\alpha) \mathbf{e}^T) \quad (3)$$

とする。

Google 関係者による示唆によれば、 $\alpha = 0.85$ が Google 推
奨値ということである¹²⁾。数学的にはこの G の最大固有値
 $\lambda_1(G) = 1$ に対応する固有ベクトル \mathbf{r} が PageRank 値となる。
これを求める数値計算アルゴリズムとしては、行列ベクトル
積計算とノルムによる正規化計算のみから成るべき乗法¹³⁾ を
用いるのが適切である。最も計算量を必要とする行列ベクト
ル積では、 H がランダム疎行列であることを考慮すると劇的
に計算性能が上がる。従って、(3) 式に基づいて行列ベクト
ル積の計算を行うことが望ましい。

3.2 URI 間の関連性と検索キーワードによる小規模な順序 変更

基本的に、Google では PageRank に基づいて検索結果の表
示順を決めていることになっているが、検索キーワードや URI
によっては局所的にそれが破られることもある。

例えば検索キーワードとして「まんねん貧乏」と「得能史
子」を与えた時の検索結果を Fig.4 に示す。

Amazon の URI と、「うだうだ Weblog」の URI の表示順が、
検索キーワードによって入れ替えられていることが分かる。厳
密に PageRank だけで表示順を決めているのであれば、局所
的とはいえ、このような順序変更は起こらない筈である。ま
た、インデントを付けることで関連が強いと判断された URI
の引っ張り上げも行われており、これは意図的に PageRank を
無視して実行されている。

しかし、このような観測を行っても表示順の詳細は不明な
ままである。我々のサーチエンジンにおいても、PageRank 以
外の要素をどのように使用してユーザの望む順に表示するか
はまだ詰め切れていない。従って、現時点ではこの機能は実
装していない。

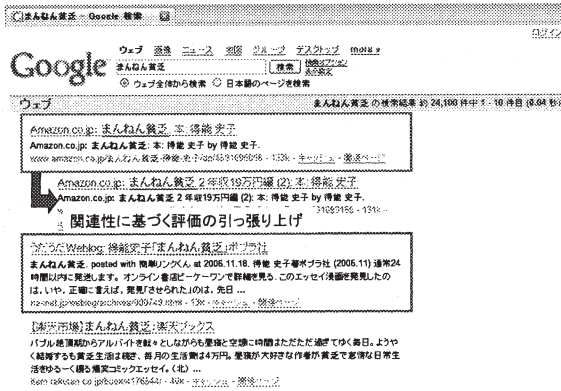
4. サーチエンジンシステムの概要

以上のような Google の機能に基づき、最終的に作成した
サーチエンジンシステムは Fig.5 のような構成になった。

教育用のシステムであるため、最小限の 3 つのコンポー
ネント構成で実現できるようにしてある。各コンポーネントは
次のような機能を持つスクリプトもしくはソフトウェアを使用
している。

DBMS フリーで使用できる MySQL¹¹⁾ を採用する。収集
した全てのデータはここに格納される。

Google検索例(1/2)



Google検索例(2/2)

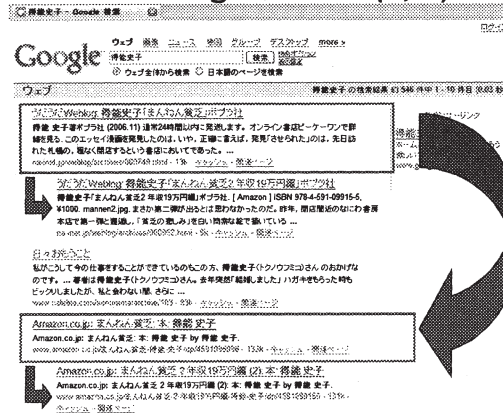


Fig. 4: Google 検索例— 検索キーワードによる順序変更 (左: 「まんねん貧乏」、右: 「得能史子」)

検索用 IF 現状ではまだ仕様が固まっていないが、最終的には PHP スクリプト (+ JavaScript?) による実装を予定している。検索結果は後述するランキングの高い順に表示されるようになる。

Web ロボット Perl による Web データ収集プログラムで、現在のところ 400 行程度の規模である。robots.txt による収集制御に対応している。

図では DBMS、検索用 IF、Web ロボットを別マシンで動作させているが、実習ではこれらを 1 台の仮想マシン上で動作させることになる。

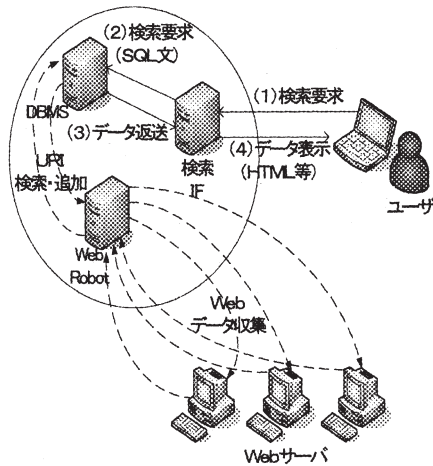


Fig. 5: サーチエンジンのシステム

これらのうち、特に重要な機能である Web ロボットと、ランキング計算部分について以下に述べる。

4.1 Web ロボット

Web ロボットの基本動作は

1. 既知の URI にアクセスし、テキストデータのみを取得

2. 取得したテキストデータ (HTML など) から新たな URI を取得 ⇒ 1 へ

というきわめて単純なものであるが、一歩間違えると悪質なアタックツールにもなりかねない危険な側面を持つ。そのため、次の機能を持つことが、安全な Web ロボットの運用には不可欠である。

- 同一 Web サーバにアクセスする際には、必ず十数秒以上の間隔をあける。
- robots.txt⁹⁾ によるアクセス制御に従う。

更に、新たな URI 取得のためには meta タグやアンカーを適切に読み取る機能が欠かせない。このためには HTML Parser が必要である。また、検索のためには Unicode に対応したエンコード機能を持ち、格納するテキストデータのエンコード方式を統一しなければならない。これから増えると思われるオープンドキュメント形式のオフィス文書ファイルも取り込むためには圧縮ファイルへの対応も求められるだろう。

現在のところ、これら全ての機能を満足する機能を備えた言語環境はそう多くない。そのため、パフォーマンスの向上に関しては難はあるものの、今のところは Perl とそのモジュール群を利用して Web ロボットを実装してある。この結果、スクリプトサイズも 500 行に満たない程度に抑えられている。

それでも学生は慣れない言語とモジュール群を活用しなければならないため、実際の PBL を行う際には配慮が必要となる。まず Perl 言語の解説を行い、次に DBI モジュールの使い方を通じて MySQL の操作手法に慣れ、その後に Encode, LWP, RobotUA, HTML Parser モジュールの機能を、短いスクリプトを作り動かすという経験を通じて理解させるようにする。慣れない言語と高度なネットワークアプリケーションの機能に親しむには、このように段階を追って少しずつ Web ロボット全体の機能が把握できるようにする配慮が必要であろう。

4.2 ランキング計算処理 — BNCsparse ライブラリの開発

前述したように、ランキングは Google 行列 G をべき乗法にかけることで得られる。我々のサーチエンジンではこのランキングを WebRank と呼ぶことにした。

はより高いものが求められる。どの要素もサーチエンジン性能に直結するものであるから、本来ならばそれらを全てチェックしなければならないが、現状では用意できる環境には制約があるので、その範囲内で出来ることを行って見たのが以下に述べる結果である。

5.1 本学内 LAN からの URI 収集と WebRank 計算結果

まず静岡理工科大学研究実験棟 LAN から URI 収集を行った結果 (Fig.7) と、それに対して WebRank 計算を行った結果を示す。

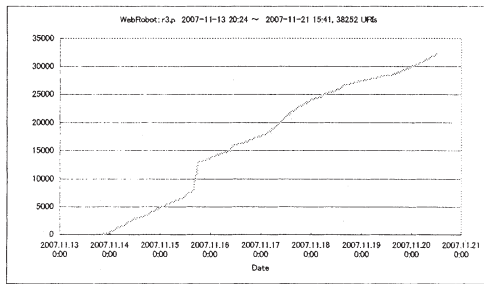


Fig. 7: 学内 LAN 内からの URI 収集数

条件は以下の通りである。

- 静岡理工科大学 LAN からアクセスしたので、Firewall と研究実習棟 Proxy の 2 段の Proxy を介して外部 7Mbps 回線を用いたことになる。
- 使用ハードウェア・ソフトウェア環境は次の通り
 - CPU Intel Pentium 4 (1.8GHz)
 - RAM 1GB
 - HDD 40GB
 - OS CentOS 5 x86_32 版
- WebRobot(r3.pl) は前述したように Perl スクリプトとして実装し、1 プロセスのみ起動。完全なシリアル動作である。
- Yahoo! Japan のトップページ (<http://www.yahoo.co.jp/>) からスタートした。後述するように、現状の WebRobot の収集アルゴリズムは深さ優先になっているため、殆どの URI は Yahoo! Japan サイト内のものとなっている。
- 1URI にアクセスするごとに 10 秒の wait を置くようにした。
- robots.txt の指示に従ってアクセス制限をしているので、意図的に WebRobot を拒否しているサイトの情報は収集できない。
- テーブルに記憶する URI は接続確認されたもののみである。アクセスできなかったものは記憶しない。

この結果、一週間で約 3 万 URI が収集できた (Fig.7)。更に収集を行い、最終的には約 8 万 URI 集めるに至った。これらのリンク構造 (行列 H) を Fig.8 に示す。黒点は非ゼロ要素があることを示している。

前述したように、3 万 URI でも 8 万 URI でも殆どが Yahoo! Japan サイト内の URI なので、行列の下三角部分には殆どリンクがない。これは発見した順に URI に固有 ID 番号を振っているためだと思われる。

これら 2 つのリンク構造に対して WebRank 計算した結果を次に示す。WebRank 計算のためのべき乗法の収束条件は次のように設定した。

$$\|r_k - r_{k-1}\|_1 \leq 10^{-5} \|r_{k-1}\|_1 + 10^{-50} \quad (4)$$

高速な CPU を持つマシン (Intel Core2Duo 2.16GHz, gcc 4.0.1 の場合) で WebRank を計算した結果は次の通りである。

1. テキストファイルからの読み込み
 - 38243URIs: 0.71 sec
 - 81305URIs: 1.78 sec
2. べき乗法の計算時間
 - 38243URIs: 0.09 sec (11 回反復, $\alpha = 0.85$)
 - 81305URIs: 0.30 sec (16 回反復, $\alpha = 0.85$)

現状では、データベースからリンク構造を引っ張り出してファイルに書き出し、それを読み取って J を構築している。従って、テキストファイルの読み込み時間とべき乗法部分の計算時間の和が WebRank 計算全体の時間を規定していることになる。この結果によって、テキストファイルからの読み込み時間はべき乗法の計算時間に比べて約 6 倍~7 倍かかっていることが分かる。よって、べき乗法部分をこれ以上高速化しても効果薄であり、それよりは I/O 性能を上げる方がよい。

5.2 URI 収集性能

更なる高速化を目指すため、URI 収集においてボトルネックになりそうな部分の改善を行って、実験を行った。改良したのは以下の 3 点である。

- ネットワークをより高速な外部回線 (ベストエフォート 100Mbps) を持つ所に変更
- WebRobot の処理をマルチプロセス化し、並列動作を行うようにした
- Cache メモリを搭載した RAID カードを用いて SATA HDD を RAID0 構成にし、ローカルストレージの性能向上を図った

このうち、マルチプロセス化による並列動作と、RAID 構成を用いた改善の概要について以下で解説を行う。

5.2.1 マルチプロセス化による並列動作

マルチプロセス化した WebRobot の動作概念図を Fig.9 に示す。

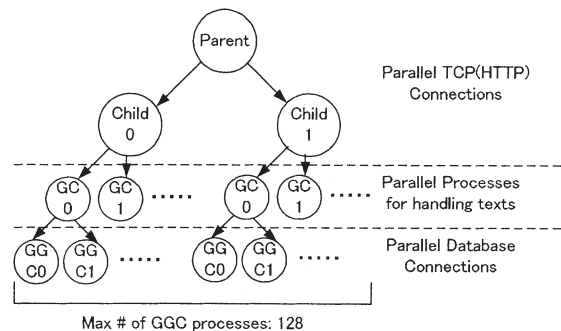


Fig. 9: Web Robot 処理の並列化

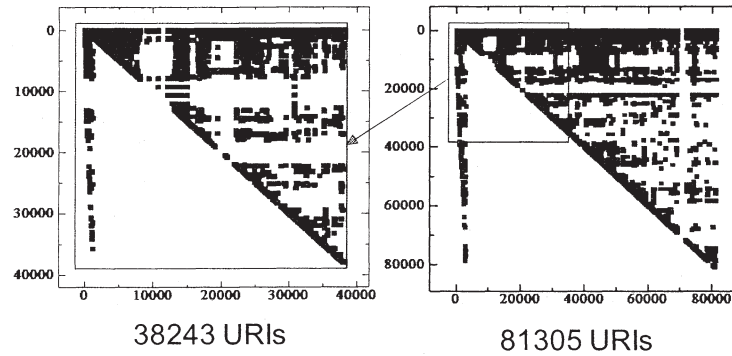


Fig. 8: H の構造

親 (Parent) プロセスは子 (Child) プロセスを fork を使って二つ起動した後はプログラム終了まで何もしない。子 (Child) プロセスは Web からのデータ収集を行うための TCP 接続 (HTTP) を担当し、接続が完了してから孫 (GC, Grandchild) プロセスを生成する。孫プロセスはひ孫 (GGC, Grand Grandchild) プロセスを生成し、テキスト処理を行った後、すぐに終了する。ひ孫プロセスはデータベースとの入出力を担当し、処理を終えた後はすぐに消えるようになっている。

全体としては親からひ孫プロセスを含めた合計数が上限値 (128) に達するまでは孫およびひ孫プロセスが生成され続ける。もし上限に達した場合は、子プロセスが 10 秒間 sleep した後、上限を下回っていることを確認できれば再度孫プロセスを fork する。

このような多段 fork によってプロセスを生成するようになったのは、子プロセス以下のプロセスをゾンビ化 (親プロセスへのリンクを失った状態) させないためである。末端のひ孫プロセスは init が預かるのでゾンビ化することはない。

一般に、近年はプロセスよりはスレッドが並列動作のために使用されることが多いが、ここではスレッドもプロセスも同じ実体として実装されている Linux を用いているので、枯れた技術である fork によるマルチプロセス構成を採用した。これによって、より大量の URI を収集することが可能になる効用が期待できる。

また、並列動作の性能向上を確実化するために、Quad-core CPU を搭載した次のようなマシン上で動作実験を行った。

- CPU Intel Core2Quad 6600
- RAM 4GB
- HDD 250GB
- OS CentOS 5 x86_64 版

5.2.2 RAID カードによるローカルストレージの高速化

これに加えて、I/O 性能の向上を図るため、大容量の Cache メモリを搭載した RAID ボードを上記のマシンに乗せ、RAID0 の機能を用いた。この結果、I/O 性能のボトルネックが解消され、より大量の URI が収集できる可能性が増える。

以上のような環境で URI 収集数の変化 (1 週間分) を Fig.10 に示す。比較のため、RAID ボードなし (No RAID0) の結果も載せてある。

RAID 機能を用いた場合は使用しない場合に比べ、最終的には 1.5 倍の URI 数の収集が可能になっていることが分かる。

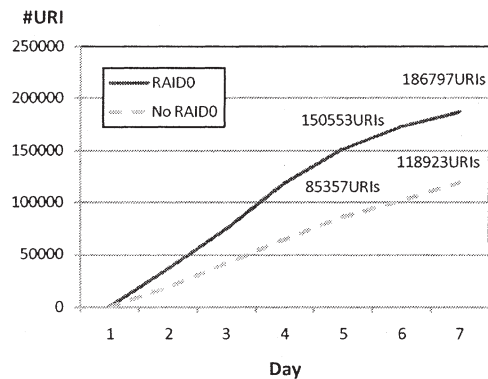


Fig. 10: 100Mbps 回線からの URI 収集数

なお、このグラフからはよく分からないが、RAID0 を使用した結果、ローカル側のルータがハングアップし、URI 数の伸びが若干削がれてしまっている。そのため、より高性能のルータを用いることにより、更なる性能向上が見込めることが予想できる。

6. まとめと今後の課題

以上、設置予定の Web デザイン特別プログラムの紹介、3 層 Web プログラミング教育の必要性とその応用としてのサーチエンジン構築の意義、そして教材として使用する予定のサーチエンジンシステムの性能評価を述べてきた。今後の課題はまず PBL の実践を行うことに尽きるわけだが、その際には以下のことをもう少し詰める必要があると我々は考えている。

教育課題 3 年生向け実験講座での教育実践の結果、PBL としても我々の目指すレベルは相当高く、学生にとっては厳しいものになると思われる。もちろん途中でドロップアウトさせないように精神的にも学習面でも注意深く指導する必要がある。そのためには、充実した内容の印刷物のテキスト²⁾ が不可欠である。これを次年度中に是非とも完成させたいと考えている。

研究課題 サーチエンジンシステム自体は、人間とインターネットとの接点となる不可欠のツールとなっており、学術研究課題としてもかなり面白いテーマであると考えている。現在のところ、次のことについて、このPBLで作成したサーチエンジンによる知見を元に研究テーマとして発展させられないか、と思案中である。

- ランキング計算の高速化・安定化と、ユーザ要求に応じたランキングの変動機能
- 日本語形態素解析を利用した検索要求部の改良

応用課題 我々が用意できるネットワークリソースを考えると、とても Google に対抗できるようなサーチエンジンシステムを作り上げることはできない。そこで、このPBLで構築したシステムを限られたネットワークエリア内の組織内ツールとして運用することを思案中である。その際には、グループウェアに取り入れられているような、組織にジャストフィットするようにランキングや Web ロボットをカスタマイズするための機構を入れる必要があるだろう。

また、外部回線さえ高速にできれば、更なる大量の URI が収集できることが判明した。多重にマシンを並べ、更なる高速化をローカルルータの限界まで行うことができれば、データマイニングのための大量のテキストデータ収集が行える。

謝辞

本研究は静岡理科大学特別研究費の支援を受けて行われた。関係者一同に厚く御礼申し上げる。

参考文献

- 1) 竹口友大・幸谷智紀, ランク機能付きサーチエンジンの開発および I/O ボトルネック対策, 第 70 回情報処理学会全国大会講演集, 2008.
- 2) 幸谷智紀・竹口友大, “サーチエンジンを作ろう”, 未公開テキスト.
- 3) M.W.Berry, M.Browne, Understanding Search Engine, SIAM, 1999.
- 4) Web デザイン特別プログラムの紹介, <http://ex-cs.sist.ac.jp/~suganuma/dep/PBL/PBL.html>
- 5) 日経データプロ編, WWW-データベース連携システム構築法, 日経 BP 社, 1996.
- 6) 幸谷智紀, VMware を用いた PC cluster と PC 教室との共存実験, 平成 18 年度情報教育研究集会論文集, pp.507 - 501, 2006.
- 7) The Community ENTerprise Operating System, <http://www.centos.org/>
- 8) Movable Type, <http://www.movabletype.com/>
- 9) <http://www.robotstxt.org/>
- 10) Google の秘密 - PageRank 徹底解説, <http://www.kusastro.kyoto-u.ac.jp/~baba/wais/pagerank.html>
- 11) MySQL, <http://www.mysql.com/>
- 12) A.N.Langville, C.D.Meyer, Google's PageRank and Beyond, Princeton University Press, 2006.
- 13) 幸谷智紀, ソフトウェアとしての数値計算, <http://na-inet.jp/nasoft/>
- 14) 山本達文, "PHP+MySQL と茶釜 (Chasen) を用いたサーチエンジンの作成", 2007 年度静岡理科大学情報システム学科卒業研究.
- 15) 伊藤裕晃, "サーチエンジンのための検索順序づけ (Rank 付け) に関する研究", 2007 年度静岡理科大学情報システム学科卒業研究.