

## 分散環境を利用したCD-ROMオンラインカタログシステム

藤田岳久 阪口哲男 杉本重雄 田畑孝一

図書館情報大学

本稿では、UNIXワークステーションによる分散システム環境の上で開発したCD-ROMオンライン目録の性能について論ずる。我々は大型計算機上のオンライン目録データベースをワークステーションに移し、それをCD-ROMと磁気ディスク双方に格納した。また、検索ソフトウェアは、サーバクライアントモデルに基づき設計・開発した。CD-ROMオンライン目録の性能を磁気ディスクに格納したものの性能と比較しながら評価した結果、CD-ROMオンライン目録は、ワークステーションのキャッシングの効果などにより、十分実用に耐えうる性能を持つという結論を得た。

### **CD-ROM Online Catalog System based on Distributed Environment**

Takehisa FUJITA, Tetsuo SAKAGUCHI, Shigeo SUGIMOTO and Koichi TABATA

University of Library and Information Science

This paper presents performance evaluation of the online catalog developed on CD-ROM called CD-OPAC. CD-OPAC is developed on a UNIX workstation and oriented to distributed environments. The authors transported the OPAC database of the University Library of ULIS from the mainframe to the UNIX workstation, and installed it both into CD-ROM and fixed disks managed by NFS. The retrieval software is configured based on the client-server model for distributed environments. Based on the experiments on CD-OPAC, this paper shows that the performance of CD-OPAC is not so bad because of caching mechanisms on the UNIX workstation and the retrieval software.

## 1. はじめに

図書館情報システムは典型的な計算機アプリケーションシステムであり、その運用には大型計算機が広く使用されている。一方、情報メディアの技術や計算機システムの性能は近年飛躍的に向上し、コストパフォーマンスが高く、強力なネットワーク機構を備えたワークステーションやパーソナルコンピュータ、大量のデータを格納できる新しいメディアなどが次々に現れてきた。そして様々な分野で、これらを利用してダウンサイジングや分散化が行われている。図書館情報システムにおいても、図書館の持つ膨大な情報を効率よく蓄積・管理し利用者に提供するために、ダウンサイジング・分散化を行うことが重要である。

本論文では、図書館情報システムのダウンサイジングの第一歩として分散システム環境の上で開発したCD-ROMオンライン目録システム(以下、CD-OPACと略す)の評価について述べる。このシステムは、UNIXワークステーション上で実現しており、ネットワークを通じて一般利用者に開放することができる。CD-ROMは大容量で持ち運びが便利な上、安価であるという特長を持ち、目録システムのような大きなデータを扱うシステムにとっては重要なメディアであり、実際にデータベースや雑誌論文目録の頒布、電子図書の出版など、商用の分野で利用されている。本論文では、CD-OPACの性能を、同じデータを磁気ディスクに格納したもの(以下NFS-OPACと略す)と比較しながら議論する。オンライン目録データは、大型計算機上で稼働している図書館情報大学附属図書館のオンライン目録システム(ULIS OPAC)[1]から転送した。

## 2. 図書館情報システムの

### ダウンサイジングとCD-ROMの利用

#### 2.1 図書館情報システムのダウンサイジング

従来の図書館情報システムは大型計算機上で開発されてきた。最近のワークステーションやパーソナルコンピュータなど個人使用向けの計算機は、柔軟なネットワーク機構を持ち、大量のデータを格納でき、高い計算能力を持っているなど、大型計算機に比べ見劣りしない性能を有している。例えば、最近のUNIXワークステーションによる分散システムで使用されているファイルサーバは、100ギガバイト以上のデータを格納することができる。

カーネギーメロン大学のProject Mercuryは、ワークステーションとパーソナルコンピュータによる分散型図書館情報システムLibrary Information System II (LIS II)を開発した[2][3][4]。LIS IIの目録システムはサーバクライアントモデルに基づき開発されており、サーバマシンでは目録データベースの管理、検索、セキュリティの保持などを行い、クライアントマシンでは検索リクエストの生成およびサーバへの送付、ユーザインタフェースの提供などを行っている。クライアントマシンは、性能に見合ったユーザインタフェースを備えたワークステーションやパーソナルコンピュータなど様々なタイプの小型計算機であり、それらは図書館本館、分館、研究室に置かれ、キャンパスネットワークを通してサーバマシンと接続されている。また、LIS IIは雑誌記事のイメージデータベースを備えている。雑誌記事はイメージリーダーで読みとってCD-ROMに格納しており、ウィンドウ上に1ページずつ表示することができる。利用者は、一つの端末から記事の検索と読書の両方を行うことができる。

Project Mercuryの実験により、LIS IIは以前使用していた大型計算機による図書館情報システムに比べてパフォーマンスが上がったという結果が出ている。また、グラフィック機能やウィンドウシステムを利用して開発したユーザインタフェースは、わかりやすく使いやすいという利用者の評判を得ているという。このことから、図書館情報システムのダウンサイジング・分散化は大変有効であると言える。

#### 2.2 CD-ROMの利用

ワークステーションなどの個人向けシステムによる分散環境においては、OPACデータベースを格納する主たるメディアとしてはファイルサーバ上の磁気ディスクがまず第一に考えられる。ファイルサーバの磁気ディスクをNFS(Network File System)のようなネットワーク指向のファイル管理システムによって管理すれば、たとえOPACが分散環境上で実現されファイルサーバが複数存在するとしても、中央集権的制御を行うことができる。

一方、分散環境においては、OPACデータベースは必ずしもファイルサーバ上の磁気ディスクに格納される必要はなく、大きな容量を持った他のランダムアクセス可能なメディアに

格納してもよいということが考えられる。我々はこのようなメディアとして、CD-ROMに着目した。OPACデータベースの格納メディアとしてCD-ROMを利用することを想定すると、頻繁にアクセスされる磁気ディスクによく見られるメディア自体の障害などを低減できるのはもちろんのこと、サーバマシンが障害を起こした際に、CD-ROMのみを他のマシンに移動するだけで目録サービスを続行できる、サービスの利用が増えてくればCD-ROMを複製してサーバを増やすことができるなど、様々なメリットが考えられる。

CD-ROMを利用したシステムにおいて問題になるのは、CDのアクセス速度の遅さである。我々は本研究において、分散環境におけるオンライン目録メディアとして、CD-ROMがどの程度利用可能なものかを実験した。実験の結果、以下のようなことがわかった。検索実験において、CD-OPACは13万件の中から8件を検索するのに約3秒、512件を検索するのに約15秒かかった。一方、NFS-OPACではそれぞれ0.2秒、1.5秒となった。これは、コンパクトディスク装置と磁気ディスクの平均アクセス時間の比が約50:1であることを考えると、悪くない結果であると思われる。

CD-OPACはサーバクライアントモデルに基づき開発し、TCP/IPプロトコルで接続されたUNIXワークステーション上で動作する。クライアントは、OSF/Motifを用いて開発したグラフィカルユーザインタフェース(GUI)を提供する。利用者はウィンドウに検索コマンドを入力し、結果を他のウィンドウで見ると、結果ウィンドウは必要なだけ残しておくことができ、中間結果を次の検索へフィードバックさせて望みの最終結果を得ることができる。

### 3. CD-ROMによるオンライン目録

#### 3.1 大型計算機からのオンライン目録データの転送

ULIS OPACは大型計算機(HITAC M660K)上で管理運用されており、約13万件の目録レコードを持っている。ULIS OPACデータベースの大きさは、約800メガバイトである。

我々はCD-OPACを開発するにあたり、大型計算機上のULIS OPACデータベースをUNIXワークステーション(SONY NEWS)へ転送した。作業は、以下の手順で行った。

(1) 大型計算機上で、目録データベース中の不必要なデータを削除し、ワークステーションで扱いやすい形式に変換する。

検索用目録を作成するには不必要と思われる図書管理用データなどを取り除き、ワークステーション上での目録データベース作成に適した形式に変換した。

(2) ファイル転送機能を用いて転送を行う。

(3) ワークステーション上で文字コードの変換を行う。

本学の大型計算機で使用されている文字セットはEBCDIK(1バイト、アルファベットなど)とKEIS(2バイト、漢字・ひらがな・カタカナなど)である。標準KEISに含まれていないいくつかの漢字は、外字として登録してある。UNIXワークステーションで使用されている文字セットはASCII(1バイト)とEUC(2バイト)であり、我々はデータベースをEBCDIK/KEISからASCII/EUCに変換した。

#### 3.2 目録データベースと検索ソフトウェア

データベース作成には、全文検索対応データベース作成ソフトウェアMediaFinder DPSI<sup>5)</sup>を利用した。作成したデータベースを磁気ディスクに格納し、さらにCD-ROM(CD-WriteOnce disk)にコピーすることによって、CD-OPACとNFS-OPACの2つを開発した。図1に目録レコー

タイトル
著者名
版表示
出版地、出版者、出版年
ページ数、大きさ
シリーズ名、巻号
注記
ISBN、価格
排架場所
図書ID
件名(サブジェクト・ヘディング)

図1 目録レコード

ドの内容を示す。データベースはインデックスファイルとテキストファイルから構成されており、データベース全体の大きさは250メガバイトとなった。

検索ソフトウェアは、サーバクライアントモデルに基づき開発した。サーバ(OPACサー

パ)とは、クライアント(OPACクライアント)からの要求を受けとり、CD-ROMまたはNFS(磁気ディスク)に対する検索操作を行う。また、クライアントは、ユーザインタフェースの提供およびサーバとの通信(TCP/IPプロトコル)を行う。クライアントは、利用したいデータベースそれぞれに対して一つずつサーバを起動する。すなわち、一つのクライアントはいくつものサーバと通信することができ、一つのサーバは一つのクライアントとのみ通信する構成となっている。クライアントとサーバは同一のワークステーション上で動作してもよいし、異なるワークステーション上で動作してもよい。一つのワークステーション上では、メモリの許す限り複数のサーバおよびクライアントを動作させることができる。サーバはクライアントからの要求に従って、データベースに対して、検索、タイトル取り出し、内容取り出しという3種類の大きな操作を行う。

- 検索操作: 与えられた検索語を含むレコードの件数を返す。
- 内容取り出し操作: 指定したレコードのすべてのフィールドの内容を返す。
- タイトル取り出し操作: 指定したレコードのタイトルのみを取り出す。内容取り出し操作に比べ高速に実行される。

なお、サーバ開発に際しては、全文検索ソフトウェア開発ツールMediaFinder SDK[6]を利用した。

図2に分散環境の構成を示す[7][8][9]。磁気ディスクに格納された目録データベースは、ネットワーク上のすべてのワークステーションから遠隔アクセスが可能である。一方、CD-ROM上の

目録データベースは、遠隔からのアクセスはせずローカルサーバからのアクセスのみである。

### 3.3 グラフィカルユーザインタフェース

グラフィカルユーザインタフェース(GUI)がキャラクタインタフェースに比べよりよい会話環境をもたらすことに疑いの余地はない。クライアントは図3に示すようなグラフィカルユーザインタフェースを持つ。クライアントは、コマンドを入力するためのウィンドウ(コマンドウィンドウ)を利用者に提供する。また、コマンドウィンドウの結果表示ボタンを押すと、目録レコードの内容を表示するためのウィンドウ(結果表示ウィンドウ)を生成する。結果表示ウィンドウは利用者が消さない限り残るので、利用者は複数の結果表示ウィンドウを見ながら新しいリクエストをコマンドウィンドウに入力する、といった操作を行うことが可能である。なお、グラフィカルユーザインタフェースはOSF/Motifを用いて開発した。

## 4. CD-ROMオンライン目録の性能評価

### 4.1 実験

本章では、CD-OPACとNFS-OPACの性能を知るために行ったレスポンスタイム測定実験について述べる。両OPACの持つ目録データベースの内容とソフトウェア構成は同じであり、格納したメディアのみが異なる。実験は、以下に述べる実験条件を変化させてレスポンスタイムを測るという方法で行った。また、レスポンスタイムのみを得るため、グラフィカルユーザインタフェースは使用していない。レスポンスタイムは、クライアント内でUNIXのシステムコールを使用して測定した。なお、表1に実験に

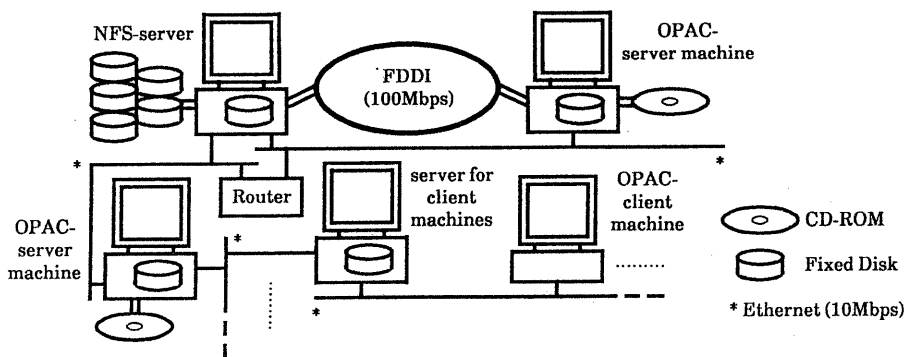


図2 分散環境の構成

表1 ワークステーションの性能特性

- サーバマシン (Sony NWS-3870)  
25MIPS, 4.3MFLOPS, 19.2SPECmark,  
128KB cache, 64MB RAM, 124MB virtual-  
memory
- クライアントマシン (Sony NWS-3410)  
17MIPS, 2.3MFLOPS, 11.9SPECmark,  
128KB cache, 12MB RAM, 37MB virtual-  
memory
- CD装置 (Sony NWP-551)  
アクセスタイム: 0.45sec (1/3 stroke),  
0.7sec (full stroke)  
データ転送レート: 150KB/sec (SCSI)
- 磁気ディスク装置 (Sony NWP-552)  
アクセスタイム: 13.5msec  
データ転送レート: 5MB/sec (SCSI)

使用したワークステーションの性能特性を示す。

以下に、変化させた条件について述べる。

(1) データベースの格納メディア

CD(CD-OPAC)と磁気ディスク(NFS-OPAC)について、実験を行った。

(2) サーバ・クライアントモデル

3.2で述べたように、一つのサーバは一つのクライアントとのみ通信する。複数のサーバが「サーバマシン」と呼ばれる1台以上のワークステーション上で稼働する。クライアントプロセスは、可能な限り別々のマシンで実行した。

実験に使用したモデルは以下の3種類に分類できる。

- 1-server-1-client: 一つのサーバをマシン上に置き、一つのクライアントを同一または異なるマシン上に置く。
- 1-server-N-clients: N個のクライアントをN台のマシン上に別々に置き、N個のサーバを一つのマシン上に置く。クライアントはそれぞれ対応するサーバと通信する。
- (M)-server-N-clients: N個のクライアントをN個のマシン上に別々に置く。M×N個のサーバプロセスをM個の目録データベースを持つ一台のマシン上に置く(図4参照)。例えば、「(2)-

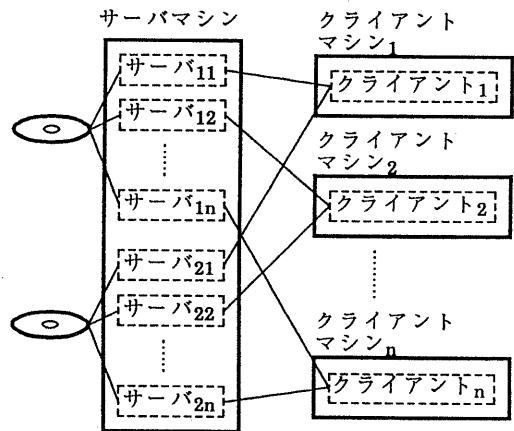


図4 (2)-server-n-clientの構成

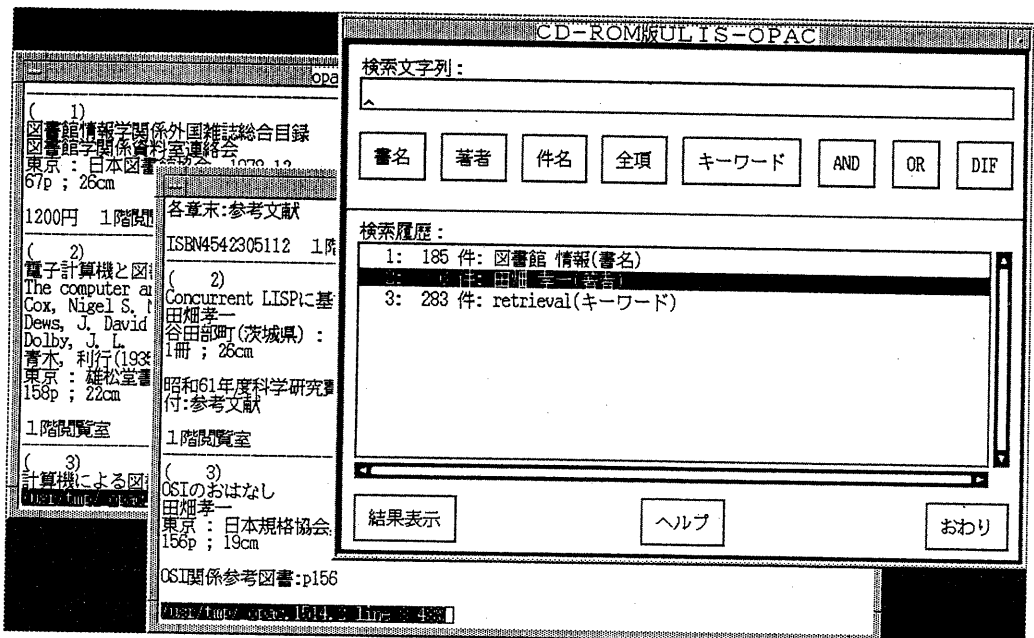


図3 グラフィカルユーザインタフェース

server-1-client」とは、2個のサーバプロセス(OPACサーバ)が2個の目録データベースを持つマシン上で動作することを意味する。なお、CD-OPACの場合、1枚のCD-ROMに1個の目録データベースが格納されているので、上記説明中「M個のデータベース」は「M台のコンパクトディスク装置」と同義である。

### (3) クライアントからサーバに対するリクエストの種類

クライアントからサーバに対する検索リクエストは、2種類準備した。これらのリクエストはサーバによって、3.2で示したデータベースに対する操作に解釈され実行される。

- 検索のみ: サーバは検索操作を行う。検索語は著者名、出版者名、書名中の語などから選んだ。サーバは検索集合に含まれるレコードの件数を返す。
- 検索およびタイトル取り出し: サーバは検索操作を行い、結果レコードに対してタイトル取り出し操作を行う。

なお、ANDやORといった集合演算操作も行うことはできるが、本実験においては使用していない。

(4) 連続するリクエストにおけるリクエスト間隔  
クライアントが2以上のサーバを持つ場合は、クライアントはすべてのサーバに対し同じリクエストを送る。「リクエスト間隔」は、利用者が検索結果を得てから次の検索リクエストを出すまでの間隔(思考時間)を表現するため、2以上のクライアントの実験において設定した。アクセス間隔は  $T-T/2$  から  $T+T/2$  の間で変化させた( $T$ はそれぞれの実験に与えられている)。

### (5) 検索語

CD-OPAC、NFS-OPACを実際に使用してみてもわかったことは、レスポンスタイムは検索結果件数が多いほど長くなる、ということである。そこで、実験には検索結果件数が8、32、128、512となる語を使用することとした。データベース中に含まれる語から検索結果件数が8件になる語を抽出し、その中からランダムに10語ずつ選んでいくつかの集合を作った。この10語から成る集合を「検索語ユニット」と呼ぶ。同様のことを、検索結果件数が32、128、512になる語について行う。実験の際には、それぞれの検索結果件数の検索語ユニットを一つずつ選び、検索語ユニットの10語を連続して(または前

項で述べたリクエスト間隔をはさんで)クライアントからサーバに送出し、個々のレスポンスタイムおよび平均レスポンスタイムを得た。

以上の5つの条件を組み合わせ、表2に示す8種類の実験を行った。

## 4.2 性能評価

前節で示した実験の結果から、CD-OPACの性能の評価を行った。

評価1: CD-OPACのレスポンスタイムから見た実用性 その1

図5は、実験8(3検索語ユニットを連続して送出)の片方のクライアントのレスポンスタイムを示したものである。15語目以降はCD-OPACとNFS-OPACにほとんど差がない。CD-OPACについて、連続して検索を行うと明らかにキャッシングの効果が現れていることがわかる。他のCDを用いた実験(実験1、2、3)においても、検索語ユニットの最初3語と残り7語のレスポンスタイムは、後者の方が短いという結果が多く見られた。

評価2: CD-OPACのレスポンスタイムから見た実用性 その2

評価1では検索およびタイトル取り出しについての評価を行ったのに対し、評価2では検索のみのレスポンスタイムを見る。図6は、実験5、6、7のCD-OPACを用いた実験について、各検索語ユニットの平均レスポンスタイムを表したものである(実験7についてはリクエスト間隔が10秒のもののみ)。これより、CDを用いた検索操作のみのリクエストの平均レスポンスタイムは2秒以下であることがわかる。最も時間がかかる検索でも、レスポンスタイムは10秒以下であり、それは、検索語ユニットの最初3語による検索の場合がほとんどである。検索操作のみのリクエストは実際のOPACシステムにおいて最も頻繁に出されるリクエストであることから、CD-OPACは実用に耐えうるシステムであると言える。

評価3: 検索リクエストの衝突について

実験7の結果を図7に示す。6個のクライアントによるCDに対するリクエストの衝突の影響がよく表されている。実際の検索においては、利用者が検索結果を得てから次の検索リクエストをキーボードから入力するのに5秒から10秒程度かかると考えられるので、グラフに示されてい

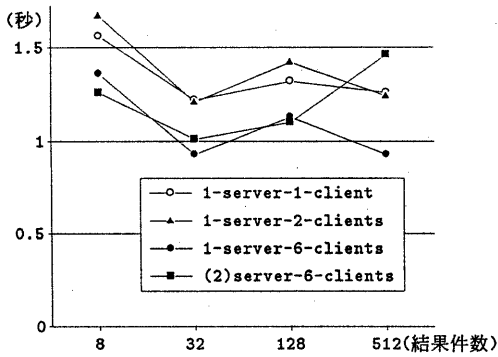


図6 CDの検索のみのレスポンスタイム

るように、あまり衝突の影響は受けないものと考えられる。なお、NFS-OPACについては衝突の影響は確認できなかった。

評価4: サーバマシンの性能について

図8は、実験1、4における各検索語ユニットの平均レスポンスタイムを表したものである。実

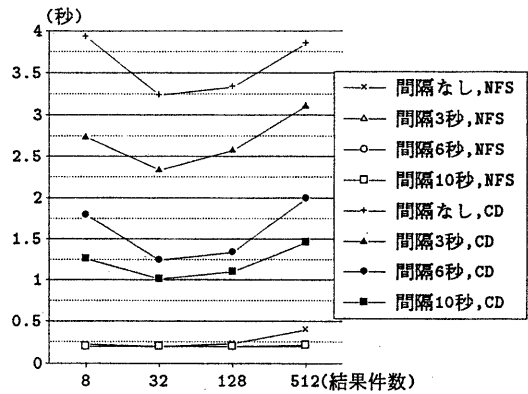


図7 思考時間の変化によるレスポンスタイムの変化((2)server-6-clients, 検索のみ)

験1と実験4は、1つのサーバマシン上で稼働するOPACサーバの数がそれぞれ1、2である。また、実験4において、2つのOPACサーバには同じリクエストが同時に送られる。グラフより、

表2 実験の条件 [実験8は、3検索語ユニット(30語)を連続して送出]

	メディア	モデル	リクエスト	間隔
実験1	CD, NFS	1-server-1-client	検索およびタイトル取り出し	
実験2	CD, NFS	1-server-2-clients	検索およびタイトル取り出し	10, 20, 30秒
実験3	CD, NFS	1-server-3-clients	検索およびタイトル取り出し	10, 20, 30秒
実験4	CD	(2)server-1-client	検索およびタイトル取り出し	
実験5	CD, NFS	1-server-1-client	検索のみ	
実験6	CD, NFS	1-server-N-clients, N=2, 6	検索のみ	10秒
実験7	CD, NFS	(2)server-6-clients	検索のみ	0, 2, 5, 10秒
実験8	CD, NFS	1-server-2-clients	検索およびタイトル取り出し	10秒

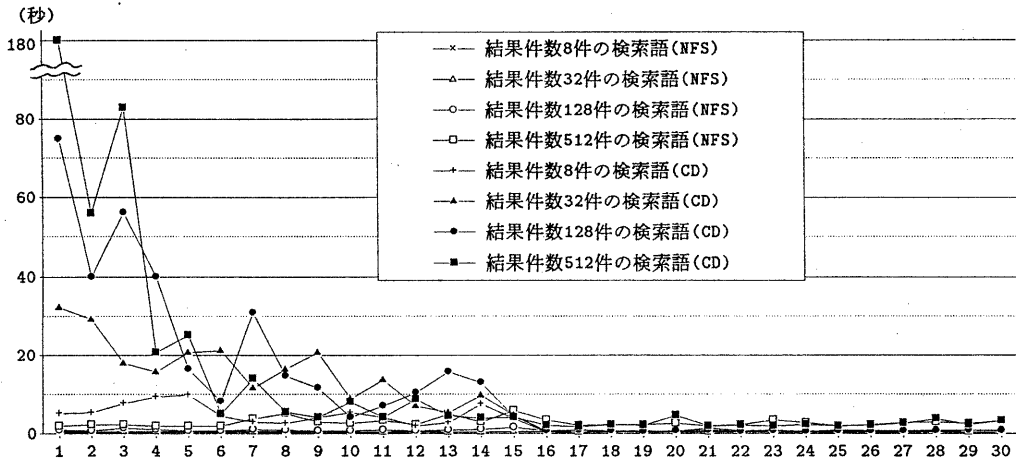


図5 30語連続検索した際の個々のレスポンスタイム

(1-server-2-clients, 検索およびタイトル取り出し, リクエスト間隔10秒)

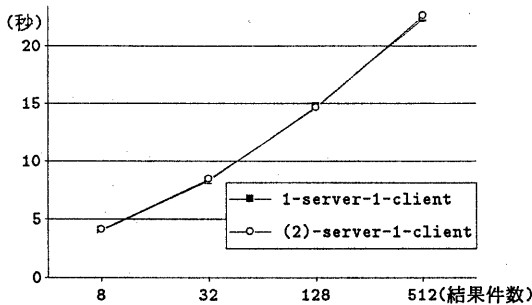


図8 1サーバマシン上のサーバ数による影響(CD, 検索およびタイトル取り出し)

実験1と実験4ではほとんど同じレスポンスタイムが得られることがわかる。サーバマシンは2個のOPACサーバを並行して実行するに十分な能力を備えている。これにより、蔵書が増加して目録がCD-ROM1枚に収まらず2枚になった場合でも、1枚の場合と変わらない性能が得られることがわかる。

## 5. 議論および結論

### [サーバの実現方式について]

今回開発した検索システムは、クライアントがデータベースを操作するためにサーバを起動するという方式をとっている。この方式では、一つのデータベースに対して複数のサーバが存在することになり、また、実験および評価でわかったキャッシングの効果を最大限に得ることができない。一つのデータベースに対し一つの常駐型サーバが動作し、複数のクライアントからのリクエストを処理する方式が望ましい。現在、常駐型サーバ方式の検索システムの設計および開発を行っている。

### [OPACデータベースをいくつかに分けることについて]

今回開発したCD-OPACでは、ULIS OPAC全体を1枚のCDに収めることができたが、蔵書が増えれば1枚に収まらなくなることも考えられる。この点については、目録データベースの各レコード間には相互参照がないので、目録データベースをいくつかに分けて別々のCDに格納し、1個のサーバマシンに割り当てることで対処できる。この場合、クライアントは複数のサーバに対して同時にリクエストを送ることができ、また、サーバ群は検索操作を並列に実行することができる。

### [OPACデータベースに新たなレコード(新着図書などの)を追加することについて]

一般的に目録は毎日更新されるものである。しかしながら、CDは1回の書き込みのみで追加のできないメディアである。我々は、蓄積されたレコードを扱うためのCD-OPACサーバと、新たに加えられるレコードを扱うためのNFS-OPACサーバ両者にアクセスできるようにOPACシステムを開発した。新たなレコードは磁気ディスク上のデータベースに追加し、クライアントからCD-OPACサーバとNFS-OPACサーバ両方を起動しリクエストを送るようになれば問題はない。

本論文において我々は、分散環境上でCD-ROMを用いて開発したOPACの性能を評価し、実用に耐えうる性能を持つことを見出した。さらに、CD-ROMによるオンライン目録は、マルチメディア指向のオンライン目録や一次情報を収める電子図書への拡張性・発展性を有している。このような点から、これからCD-ROMが図書館情報システムのためのメディアとして重要なものになり、受け入れられるであろうという実感を果たした。

### 参考文献・参考文献

- [1] 加藤信哉. 新しいオンライン目録—ULIS-MARCからULIS OPACへ. 図書館情報大学附属図書館報. Vol. 7, No. 4, p. 8-9(1991)
- [2] Troll D.A. Library Information System II Progress Report and Technical Plan. CMU, 1990, 44p. (Mercury Technical Reports Series, No.3)
- [3] The Mercury team. The Mercury Electronic Library and Library Information System II The First Three Years. CMU, 1992, 20p. (Mercury Technical Reports Series, No. 6)
- [4] 安達淳, 橋爪宏達. 欧米における「電子図書館」プロジェクト. 情報処理. Vol. 33, No. 10, p. 1154-1161(1992)
- [5] MediaFinder DPS使用説明書. 東京, ソニー, 1990.
- [6] MediaFinder SDK使用説明書. 東京, ソニー, 1990.
- [7] 杉本重雄, 田畑孝一. マルチメディアネットワークシステム: マルチメディアを活用した新しい教育と研究の環境. ビジネス・コミュニケーション. Vol. 28, No. 10, p. 28-33(1991)
- [8] 藤田岳久, 阪口哲男, 松本紳, 杉本重雄, 田畑孝一. 図書館情報大学におけるマルチメディアを指向した新しい教育と研究の環境. 情報処理学会研究報告. Vol. 92, No. 52, p. 1-8(1992)
- [9] 藤田岳久, 阪口哲男, 松本紳, 杉本重雄, 田畑孝一. マルチメディアを指向した新しい教育と研究の環境. 図書館情報大学研究報告. Vol. 11, No. 2. (発表予定)