

平成 15 年 2 月 21 日

未踏ソフトウェア創造事業
報告書

「Web 上のデータを中心とした
複数論文データベースの統合」

難波 英嗣 (広島市立大学 情報科学部)

奥村 学 (東京工業大学 精密工学研究所)

齋藤 豪 (東京工業大学 精密工学研究所)

目 次

1. 背景及び目的	5
2. 概要.....	5
3. 開発内容	7
3.1. システムの対象	7
3.2. 性能	7
3.3. 機能詳細	8
● 管理画面認証機能	8
● 論文収集機能.....	9
● データソース設定機能	9
● 本文データ管理機能.....	9
● 書誌情報・未解決参照情報抽出機能	9
● 書誌情報・未解決参照情報収集機能	9
● 参照データ生成機能.....	9
● 書誌情報管理機能	9
● 参照データ管理機能.....	9
● 論文検索機能.....	10
● 参照関係表示機能	10
3.4. データ構成	10
● 論文ファイル.....	10
● 本文データ	10
● 書誌情報.....	10
● 未解決参照情報	11
● 参照関係データ	11
● データソース設定	11
3.5. システムの構成	11
3.5.1. システムの構成.....	11
3.5.2. サーバ構成	12
● PRESRI サーバ.....	13
● クライアント.....	13
● 論文 DB サイト	13
● インターネット上の論文サイト.....	13

3.6. モジュール構成	13
● ロボット	13
● 書誌情報抽出エンジン	14
● 参照情報作成エンジン	14
● マスタ書誌情報修正 UI	14
● サーバ管理 UI	14
● データベース	15
● 検索サーバ	15
3.7. ハードウェア・ソフトウェア構成	15
3.7.1. プラットフォーム	15
3.7.2. ソフトウェア・ツール	15
3.7.3. クライアント環境	16
3.8. システムの制限	16
3.8.1. 認証機構について	16
3.8.2. 論文本文へのリンクについて	16
3.8.3. アブストラクト及び参照部分近辺の引用の表示について	17
3.9. 画面の構成	17
3.9.1. 画面の詳細	19
3.9.2. 検索画面	22
3.10. 開発計画	23
3.10.1. 開発環境	23
3.10.2. ツール	23
3.10.3. 開発の体制	23
● 書誌情報抽出エンジン	24
● 参照情報作成エンジン	24
● マスタ書誌情報修正 UI	24
● サーバ管理 UI	24
● データベース	24
● 検索サーバ	24
3.11. テスト	25
3.12. 評価	27
4. 実施計画書内容との相違点	27
5. 開発成果の特徴	28
6. 今後の課題・展望	29

7. 開発分担	29
8. 付録.....	30
9. 参考文献	41

1. 背景及び目的

特定分野の研究動向を知るためには、その分野の論文を網羅的に収集する必要がある。このような文献調査を行うのに、しばしば論文データベースが利用される。しかし、論文データベースが分散して存在していると、データベース毎に検索するのは非効率的である。そこで、本プロジェクトでは、複数の論文データベースを統合的に利用検索できるシステムの開発を行う。

また、統合システムを用いて効率的にサーベイ支援(検索)を行うためには、論文の提示方法についても検討する必要がある。ResearchIndex, Cora をはじめとする多くの引用文献データベースの論文検索インタフェースは、検索結果や参照関係にある論文をリスト形式で表示するのが一般的であったが、このような表示方法では、より大きな参照構造の中での個々の論文の位置付け(関係)がわかりにくいという問題点があった。本プロジェクトでは、論文間の参照関係をグラフで表示し、ユーザがグラフ上の論文アイコンにカーソルを重ねるとその論文の情報が、参照関係を示す矢印にカーソルを重ねると参照個所が提示できるようにする。

本報告書では、まず 2 節で概要について述べる。次に 3 節では開発内容を、4 節では開発成果の特徴を述べる。5 節では今後の課題と展望について述べる、6 節では実施計画書内容との相違点について述べる。開発分担については 7 節でふれる。

2. 概要

我々は、これまでに Web 上に存在する Postscript および PDF 形式の日英論文データを収集して論文データベースを構築している。しかし、研究者が利用可能な論文データベースは、このような Web 上の論文データ以外にも数多く存在する。例えば、近年では、国際会議や学会の全国大会では予稿集の代わりに CD-ROM が配付されることが多い。また、研究者は、それぞれの所属する組織の図書館にある論文データベース等を利用することができる。この他に、学会や出版社の所有する論文データベースも利用できる。このようなローカルに存在する論文データを、PRESRI と統合的に利用できれば、非常に便利である。例えば、CD-ROM 中の論文と PRESRI 中の論文が参照関係にあれば、その参照関係をたどって、効率的に関連論文を集めることができる。以後、個人の所有する CD-ROM、図書館や出版社や学会が所有するデータベースをまとめて、ローカルな論文データベースと呼ぶ。本プロジェクトでは、このようなローカルな論文データベースと PRESRI とを統合的に利用できるようにシステムの構築を行う。

PRESRI には、論文毎に(1)論文表題・著者名等の書誌情報、(2)インターネット上の本文データ(Postscript or PDF ファイル)の所在(URL)、(3)概要の3種類の情報が保持されている。また、(4)論文間の参照関係と参照タイプに関するデータも保持されている。PRESRI とローカルな論文データとの統合システムを構築するには、まずローカルな論文データから同様に(1)~(4)のデータを抽出し、次にローカルデータベースから抽出した情報を、PRESRI 上の論文データと統合する。その際、本プロジェクトでは2種類の統合方法を検討する。一つはローカルな計算機上でデータを統合する方法であり、もう一つは PRESRI 上でデータを統合する方法である。統合にこのような2つの方法をとる理由は、ローカルな論文データベースの著作権等の問題に関連する。例えば、ユーザが CD-ROM の論文データベースを所有している場合を考える。このようなデータの多くは、所有者が個人で利用することが前提になっているため、このデータベースから抽出した情報は、不特定多数の研究者が利用可能である PRESRI 上に置くことが出来ない。このような性質を持つ論文データベースは、ローカルな計算機上でデータを統合する必要がある。一方、論文データベースの作成者から、PRESRI 上での(部分的な)利用の許諾が得られた場合、これらのデータは PRESRI サーバ上で統合できる。

統合システムを用いて効率的にサーベイ支援(検索)を行うためには、論文の提示方法についても検討する必要がある。ResearchIndex, Cora をはじめとする多くの引用文献データベースの論文検索インタフェースは、検索結果や参照関係にある論文をリスト形式で表示するのが一般的であったが、このような表示方法では、より大きな参照構造の中での個々の論文の位置付け(関係)がわかりにくいという問題点があった。本プロジェクトでは、論文間の参照関係をグラフで表示し、ユーザがグラフ上の論文アイコンにカーソルを重ねるとその論文の情報が、参照関係を示す矢印にカーソルを重ねると参照個所が提示できるようにする。また、ユーザの用途に応じて、参照構造をズームインやズームアウト出来る機能も提供する。さらに、ユーザがグラフ上の論文アイコンにカーソルを停止させた時間と、その論文に対するユーザの関心の強さに相関があると考え、カーソルの停止時間に応じて、表示する論文情報の詳細度を変える機能を提供する。

本プロジェクトの当面の運用形態の目標は、個人の所有する CD-ROM 等の論文データ、東京工業大学図書館が所有する論文データベースと PRESRI との統合であるが、将来的には、より多くの論文データベースと統合し、最終的には世界中の論文誌を網羅した学術ポータルをユーザに提供するような運用を目指していきたい。このような環境をユーザに提供することで、効率的でかつ網羅的なサーベイが可能になり、それが結果的に研究促進にもつながると考えられる。

3. 開発内容

3.1. システムの対象

本システムは利用者及び管理者として下記を想定している。

- システム管理者

管理者は、システムのインストール及びデータの保守を行うユーザである。管理者は下記の能力を有する必要がある。

- PRESRI サーバをインストールするハードウェア及び OS について運用及び管理の基礎知識を有する
- PRESRI サーバで利用する諸ツールのインストール及び設定方法について基礎知識を有する
- PRESRI サーバ上に諸ツール等をインストールする権限を有する
- PRESRI サーバ上に PRESRI サーバが利用するデータ類の保存領域を確保する権限を有する
- Web サーバの管理についての知識を有する
- TCP/IP 及び HTTP に関する知識を有する

- システム利用者

利用者とは、本システムを利用して論文及び参照関係の調査を行うユーザである。利用者は下記の能力を有する必要がある。

- Web ブラウザを操作することができる
- 論文の構成等に関する基礎知識があり、用語が理解できる

3.2. 性能

- サイジング

本システムは将来的に以下のボリュームのデータを取り扱うことを想定している。

本文:	最大 100 万件
書誌情報:	最大 1500 万件
参照情報:	最大 2500 万件

本文データを最大 100 万件程度と想定しているのは、以下に述べる根拠に基づいている。NEC が開発・公開している ResearchIndex は、Web 上の PostScript および PDF 形式の

論文を 1 年前の段階で約 44 万件収集していることが分かっている(現在, 何万件収集しているのかは分からなかった)。このデータベースでは英語論文のみを対象にしているが, 本プロジェクトでは, さらに日本語をはじめとする多くの言語で記述された論文も対象に考えているため, Web 上のフルテキストデータは 44 万件以上になると考えられる。また, PRESRI は他の多くのデータベースと統合することを考え, 100 万件程度のデータが扱えるシステムの構築を目指す。

● パフォーマンス

本システムは以下のパフォーマンスを発揮する。

一日あたりのヒット数:	2 万件
最大同時アクセス数:	10 件
ワーストケースのレスポンス:	15 秒

単一のサーバでの目標達成が困難な場合にサーバを複数台構成にする可能性も念頭においてシステムの詳細設計及び実装を行う。

本システムでは下記の機能を提供する。

1. 管理画面認証機能
2. 論文収集機能
3. データソース設定機能
4. 本文データ管理機能
5. 書誌情報・未解決参照情報抽出機能
6. 書誌情報・未解決参照情報収集機能
7. 参照データ生成機能
8. 書誌情報管理機能
9. 参照データ管理機能
10. 論文検索機能
11. 参照関係表示機能

次節で各機能について概説する。

3.3. 機能詳細

本節では本システムが提供する各機能について概説する。

● 管理画面認証機能

管理権限を有しない利用者がシステムの設定を変更できないようにする機能。アカウント

名とパスワードによる認証を基本とする。この機能には、システム管理者の登録、削除等の機能も含まれる。

- **論文収集機能**

インターネット上で公開されている PDF, PS 等の論文ファイルを収集する。検索エンジンなどを利用して公開論文の URL リストを取得し、その結果に基づいて論文ファイルを収集する。

- **データソース設定機能**

大量の論文を保持しているネットワーク上のサーバ及び、ローカルのファイルシステムに保存されている論文集等のデータソース設定を管理する。データソース設定の作成、変更、削除機能を提供する。また、データソースごとの書誌情報更新履歴等も管理する。

- **本文データ管理機能**

ファイルシステム上に保存されている論文ファイルから本文データのテキストデータを抽出し、保存する。

- **書誌情報・未解決参照情報抽出機能**

本文データを解析し、論文の書誌情報及び論文中で参照されている被参照論文の書誌情報一覧を出力する。

- **書誌情報・未解決参照情報収集機能**

データソースサイトで公開されている、書誌情報・未解決参照情報及び、ローカルに保存されているデータソースの書誌情報・未解決参照情報を収集し、取り纏める。

- **参照データ生成機能**

各データソースから取得した書誌情報と未解決参照情報から、被参照論文の同定処理を行い、参照データを作成する。

- **書誌情報管理機能**

各データソースから取得した書誌情報、未解決参照情報から得られた書誌情報を一括して管理する。

- **参照データ管理機能**

サーバが保持している参照データを管理する。

- **論文検索機能**

キーワード等から，該当する論文を検索する．

- **参照関係表示機能**

論文間の参照関係を分析し，表示する．

3.4. データ構成

本システムでは下記のデータを取り扱う．

1. 論文ファイル
2. 本文データ
3. 書誌情報
4. 未解決参照情報
5. 参照関係データ
6. データソース設定

以下に，それぞれのデータの詳細を概説する．

- **論文ファイル**

論文ファイルは Web など公開されている論文そのものである．本システムでは PDF もしくは PostScript で記述されているものを想定している．論文ファイルはクライアントに論文本文を提示する際や，後述の本文データを抽出する際に利用する．

- **本文データ**

本文データは，論文ファイル中に記述されている内容を抽出し，テキストデータに変換したものである．論文 1 編につき 1 つずつ作成され，通常はテキストファイルとして保存されている．本文データは書誌情報や参照情報の抽出の際に利用される．また，アブストラクトや引用部分近辺の文書も本文データから抽出される．

- **書誌情報**

書誌情報は，論文のタイトル，著者，所属などの情報を取り纏めたものである．また，論文ファイルの所在情報も含んでいる．このデータは本文データ 1 編につき 1 つずつ抽出され，データベースに登録される．このデータは主に論文検索や時の際に論文ファイルへのリンク表示の際に利用される．このデータはデータソースからローカルサーバに対して公開する．

- **未解決参照情報**

未解決参照情報は、ある論文（参照論文）で参照している論文（被参照論文）のタイトル、著者、掲載誌などの情報である。参照論文中の参考文献一覧から抽出される。参考文献ごとに1つずつ作成されるので、論文1編につき参考文献の本数分だけ生成される。このデータはデータソースからローカルサーバに対して公開する。

- **参照関係データ**

参照関係データは、参照論文と被参照論文との対応関係である。未解決参照情報のタイトル、著者などと、システムに登録されている論文の書誌情報とを比較し、同定処理を行う。その結果、未解決参照データと被参照論文との対応が取れたものについて、参照論文と被参照論文のIDの対応関係をデータベースに登録する。このデータは、検索結果から参照関係ツリーを表示する際に利用される。

- **データソース設定**

データソース設定は、ローカル及びインターネットのデータソースに関する設定である。ローカルデータソースに関しては論文ファイル集合の保存場所、インターネットデータソースでは、書誌情報ファイルのURL定義などを保持している。

3.5. システムの構成

3.5.1. システムの構成

本システムの構成図を図1: システム構成図に示す。図中の網掛け部分が今回の実装範囲である。

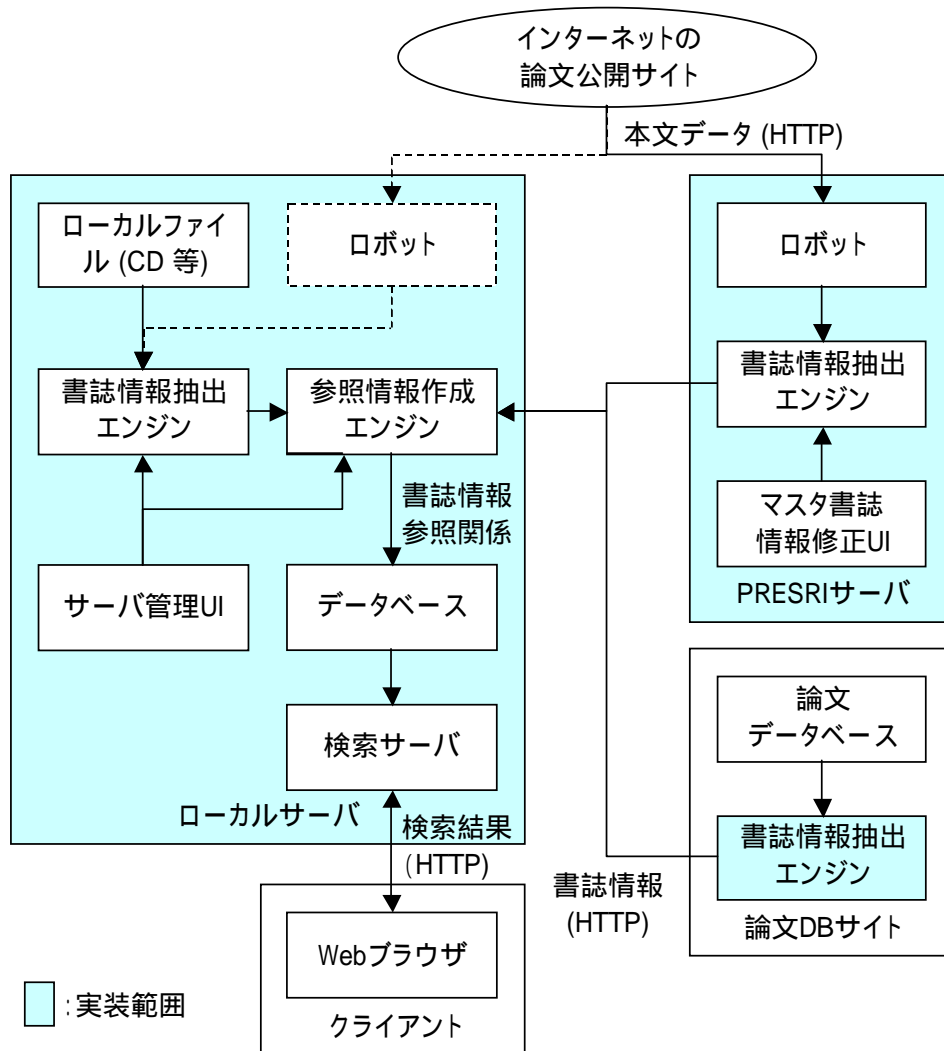


図1: システム構成図

3.5.2. サーバ構成

本システムは主に下記のサーバから構成される。

- ローカルサーバ

- 利用者に対して検索機能を提供するサーバである。
- 大学，研究機関，組織などで別個にサイトを構築し，その組織に所属する利用者にサービスを提供することを想定している。
- 複数のインターネットデータソースサイトと連携し，外部のサイトが公開する論文を検索することが可能である。
- CD-ROM などで提供される論文集データをローカルデータソースとして登録し，検索対象とすることが可能である。
- 利用者からの検索要求をうけて，検索結果や参照情報を送信する。

- 負荷分散などの目的のため、ローカルサーバを複数台のコンピュータで構成する可能性もある。
- ローカルサーバの台数構成などについては個別のサイト管理者にゆだねる。

● PRESRI サーバ

- 特定のサイトで動作する論文収集エンジンである。
- ロボット機能を利用して、インターネット上で公開されている論文を収集し、その情報をまとめて書誌情報として公開する。
- 論文を収集する際にインターネット上の検索サイトなどを利用する可能性がある。

● クライアント

- 一般的な Web ブラウザが動作するコンピュータである。
- 利用者が検索機能を利用する際に操作する。
- またシステム管理者が保守、管理作業を行う際にも利用する。

● 論文 DB サイト

- 学会の論文サイトや論文データベースサイトなどである。
- 大量の論文を保持し公開しているサイトを想定している。
- ロボットなどが論文を収集することを許容しないサイトである。
- データソースサイトとしてローカルサーバと連携動作する。
- ローカルサーバに書誌情報を提供することで、当該サイトが公開している論文書誌情報をローカルサーバから検索することが可能である。

● インターネット上の論文サイト

- 著者などが自著の論文などをインターネットで公開しているサイトである。
- 少数の論文を公開しているサイトを想定している。
- 本システムは論文公開サイトから取得した論文を検索対象とすることができる。

3.6. モジュール構成

本システムを構成するモジュールごとに機能概要を説明する。

● ロボット

本モジュールは主に論文収集機能を提供する。

検索エンジンなどを利用してインターネット上で公開されている論文を検索する。検索結果に基づいて、論文ファイルを収集する。このモジュールは原則として PRESRI サーバで

動作するが、ローカルサーバで使用することも可能である。

● 書誌情報抽出エンジン

本モジュールは主に本文データ管理機能及び書誌情報・未解決参照情報抽出機能を提供する。ローカルに保存されている論文ファイルから本文データを抽出し、指定されたディレクトリに保存する。また、抽出した本文データから書誌情報及び未解決参照情報を抽出する。

本モジュールは PRESRI サーバ、論文 DB サイト、ローカルサーバのそれぞれで動作する。PRESRI サーバではロボットが収集してきた論文ファイル集合から書誌情報と未解決参照情報を抽出する。論文 DB サイトではそのサイトが保有する論文ファイルを対象として、書誌情報等を抽出する。PRESRI サーバ及び論文 DB サイトでは抽出した書誌情報及び未解決参照情報をファイルに取り纏め、ローカルサーバに対して公開する。ローカルサーバでは、CD-ROM などの論文集から取得した論文ファイルを対象として、書誌情報等の抽出を行う。取得した書誌情報等はローカルのみで使用する。

● 参照情報作成エンジン

本モジュールは書誌情報・未解決参照情報収集機能及び参照データ生成機能を提供する。ローカルサーバに登録されているデータソース設定に基づいて、各データソースより書誌情報及び未解決参照情報を収集し、ローカルのデータベースに登録する。また、ローカルの DB 上に登録された参照情報について、未解決の被参照論文情報と DB に登録されている書誌情報とを比較し、該当する論文が見つかった場合は該当論文の ID を DB に記録する。

● マスタ書誌情報修正 UI

本モジュールは書誌情報管理機能を提供する。論文 DB サイトや PRESRI サーバがローカルサーバに対して公開している書誌情報の抽出エラーなどを修正するインタフェースである。

● サーバ管理 UI

本モジュールは管理画面認証機能、データソース設定機能及び書誌情報管理機能を提供し、ローカルサーバ上で動作する。ローカルサーバ上のデータソース設定や、各データソースからの書誌情報取得の設定、参照情報作成の設定などを変更することが可能である。また、ローカルデータソースからのデータコピーや書誌情報抽出、インターネットデータソースからの書誌情報取得、参照情報の更新を直接実行できる。ローカルに登録されている論文書誌情報について自動抽出のエラーなどを修正する機能も提供する。

データソース管理画面はユーザ認証によって保護されており、管理者以外は設定の変更ができない。データソース管理画面には管理者情報の管理機能も含まれる。

● データベース

本モジュールは管理画面認証機能，データソース設定機能，書誌情報管理機能及び参照情報管理機能を提供する．また，論文検索機能及び参照データ生成機能の一部も提供する．このモジュールはローカルサーバ上で動作する．ローカルサーバが保持する各データソースの設定や管理者の設定情報を保持している．また，参照情報作成エンジンが収集した書誌情報及び，解決済みの参照関係データを登録する．参照関係の解決や論文検索の際には，論文の書誌情報等をキーとして該当する論文のリストを検索する．また，参照関係の表示等の際には論文 ID をキーとし，参照論文及び被参照を検索する．

● 検索サーバ

本モジュールは，論文検索機能及び参照関係表示機能を提供する．主にローカルサーバ上で動作する．利用者が Web ブラウザを利用して送信した検索要求に対して，該当する論文一覧を送信する．また，論文参照関係の表示依頼に対して，論文の参照データを整形し表示データを送信する．

3.7. ハードウェア・ソフトウェア構成

本節では本システムが想定するプラットフォーム等について概説する．

3.7.1. プラットフォーム

本システムでは下記のプラットフォームを対象として想定する．

アーキテクチャ	OS	ミドルウェア等	備考
PC-AT 互換機	Linux	なし	
	FreeBSD	Linux thread	動作保証なし
	Solaris	なし	
	Windows2000	Cygwin	
	WindowsXP		
Sparc	Solaris	なし	

3.7.2. ソフトウェア・ツール

本システムは下記のツール・ライブラリを必要とする．

- リレーショナルデータベース
 - MySQL もしくは PostgreSQL
- Web ドキュメント収集
 - wget

- テキスト変換
 - prescript (PDF, PS text)
 - imdkcv (日本語コード Unicode)
- Web サーバ
 - Tomcat
 - apache
- プラットフォーム
 - sun jre
 - perl

3.7.3. クライアント環境

本システムは管理 UI，検索 UI の表示用ブラウザとして下記のことを想定する．

Web ブラウザ	バージョン等	備考
Netscape Navigator	4.7.*, 6.0, 6.2, 7.0	
Internet Explorer	5.5SP2, 6.0	
Opera	6.0	

3.8. システムの制限

3.8.1. 認証機構について

本システムでは認証を必要とする諸学会，論文データベース等へのアクセスの認証代行機能は提供しない．従って，これらのサイトから提供される情報を閲覧するためには，別途，各サイトへの認証を利用者が行う必要がある．

3.8.2. 論文本文へのリンクについて

論文検索結果などに出力される論文本文へのリンクが直接論文本文を参照できない可能性がある．具体的には下記のようなケースが想定される．

1. 論文検索結果に対する相対的な ID から論文本文を閲覧するサイト
 - 論文本文への直接のリンクを出力することが困難なため，論文検索結果画面へのリンクになる．
 - この場合，論文の検索結果へのリンクが出力される．
 - NACSIS-ELS などが該当する
2. 認証を経由せずに論文本文 URL をポイントするとエラーになるサイト
 - 論文本文へのリンクを出力しても論文を参照できない．
 - 事前にユーザが該当サイトへの認証を行っておく必要がある．
 - この場合，該当サイトの認証画面などへのリンクを出力する．
 - Europhysics Letters などが該当する

3. その他，論文本文を表示するために認証が必要なサイト

- 論文本文へのリンクを辿るとそのサイトの認証画面が表示される．
- 認証を行うことで論文本文の表示は可能である．
- 多くの学会サイトが該当する．

3.8.3. アブストラクト及び参照部分近辺の引用の表示について

論文の本文データ取得に認証が必要な論文については，論文検索結果などの表示の際に参照部分近辺の引用などを表示することができない．また，アブストラクトの取得に認証が必要な場合は，アブストラクトについても同様に表示することが出来ない．

3.9. 画面の構成

本システムの管理画面の画面遷移を図 2: 管理画面遷移図，検索画面の画面遷移を図 3: 検索画面遷移図 に示す．

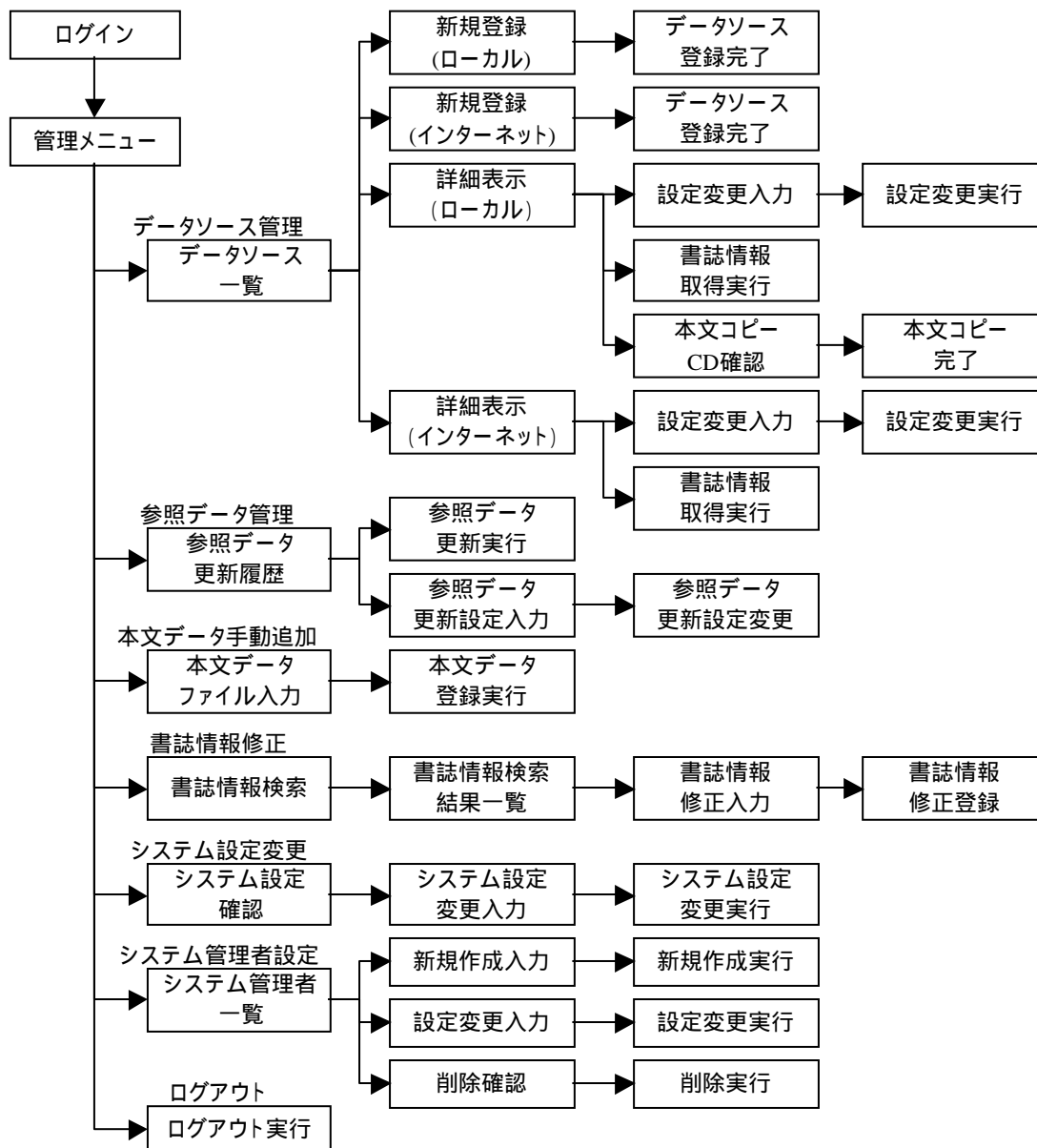


図2: 管理画面遷移図

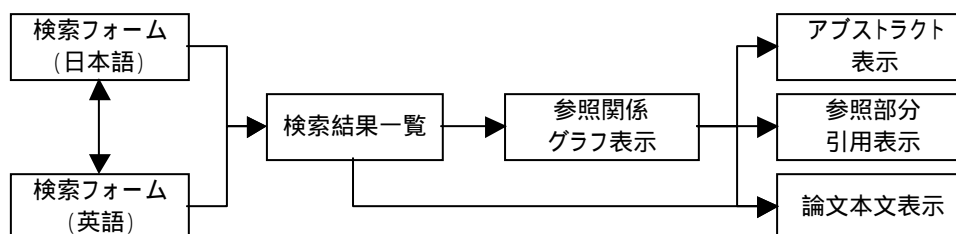


図3: 検索画面遷移図

3.9.1. 画面の詳細

各画面の機能について概説する。

3.9.1.1. 管理画面

- ログイン

管理画面のログインフォーム。

- 管理メニュー

管理画面のメニューである。メニューとして下記の項目が表示される。

- データソース管理
- 参照データ管理
- 本文データ手動追加
- 書誌情報修正
- システム設定変更
- システム管理者設定
- ログアウト

- データソース一覧

現在システムに登録されているデータソースの一覧を出力する。管理メニューからデータソース管理を選択すると表示される。

- データソース新規登録（ローカル）

ローカルデータソースを新規に作成する。本文ファイルの保存ディレクトリ等を設定する。

- データソース新規登録（インターネット）

インターネットデータソースを新規に作成する。書誌情報・未解決参照情報の公開 URL 等を設定する。

- データソース登録完了

データソース設定をシステムに登録する。

- データソース詳細表示（ローカル）

ローカルデータソースの現在の設定を表示する。

- **データソース詳細表示 (インターネット)**
インターネットデータソースの現在の設定を表示する。
- **データソース設定変更**
データソースの現在の設定を変更する。ローカル、インターネットに応じた設定項目のフォームが表示される。
- **データソース設定変更実行**
データソースの現在の設定変更をシステムに登録する。
- **書誌情報取得実行**
データソースから最新の書誌情報を取得する。
- **本文コピーCD 確認**
ローカルデータソースの論文ファイルを再取得する前に CD をマウントするよう促す。
- **本文コピー実行**
ローカルデータソースの論文ファイルをローカルのディスクにコピーする。
- **参照データ更新履歴**
直前の参照データの更新履歴及びそれ以降のデータソースの更新履歴を出力する。管理メニューから参照データ管理を選択すると表示される。
- **参照データ更新実行**
参照データの更新を即時実行する。
- **参照データ更新設定入力**
参照データの定時更新の設定フォーム。
- **参照データ更新設定実行**
参照データの定時更新設定を登録する。
- **本文データファイル入力**
本文データを手動でローカルデータソースに追加する。追加ファイルはローカルのパス、もしくは URL で指定する。指定された URL が HTML の場合は、その中のリン

クをたどる．この画面は管理メニューから本文データ手動追加を選択すると表示される．

- **本文データファイル追加**

入力された本文ファイルを取得し，ローカルデータソースに追加する．

- **書誌情報検索**

書誌情報修正を行う論文を検索するフォーム．管理メニューから書誌情報修正を選択すると表示する．

- **書誌情報検索結果**

論文の検索結果一覧．論文を選択すると書誌情報修正入力画面が表示される．

- **書誌情報修正入力**

書誌情報を修正するフォーム．タイトル等の情報を入力できる．

- **書誌情報修正登録**

修正入力画面に入力した書誌情報を登録する．

- **システム設定確認**

システムの設定を変更する．書誌情報の保存ディレクトリなどが表示される．

- **システム設定変更入力**

システム設定を変更する．設定項目の入力フォームが表示される．

- **システム設定変更実行**

システム設定変更入力画面に入力した設定を反映させる．

- **システム管理者一覧**

システム管理者を一覧で表示する．管理メニューからシステム管理者設定を選択すると表示される．

- **新規作成入力**

システム管理者を新規に追加する．アカウント，ユーザ名，パスワード，表示言語等を選択する．

- **新規作成実行**
新規作成入力画面に入力されたユーザをシステムに登録する。
- **設定変更入力**
システム管理者の個別の設定を変更する。ユーザ名，パスワード，表示言語等を選択できる。アカウントは変更できない。
- **設定変更実行**
設定変更入力画面に入力されたユーザ情報をシステムに登録する。
- **削除確認**
システム管理者の削除を確認する。
- **参照データ更新実行**
指定された管理者を削除する。
- **ログアウト実行**
管理画面からログアウトする。

3.9.2. 検索画面

- **検索フォーム**
タイトル，著者，掲載誌などの検索キーワード入力フォーム。英語と日本語のページがある。
- **検索結果一覧**
検索フォームに入力された条件に合致する論文を一覧で表示する。
- **参照関係グラフ表示**
論文間の参照関係をグラフ形式で表示する。個々の論文を年代，類似度に応じて点で表し，その間の参照関係を矢印で表示する。表示範囲については何段階かの切替が可能である。論文や矢印をポイントすると書誌情報等も表示される。
- **アブストラクト表示**
論文のアブストラクトなど，より詳細な情報を表示する。参照関係グラフからポップアップ的に表示されることを想定する。

- 参照部分引用表示

論文の参照部分近辺を引用して表示する。参照関係グラフからポップアップ的に表示されることを想定する。

- 論文本文表示

該当する論文ファイルを表示する。ローカルデータソースの論文については直接表示する。また、インターネットデータソースの論文は、論文本文ファイルへ至るリンクを出力する。

3.10. 開発計画

3.10.1. 開発環境

本システムは下記の環境にて開発を行う。

- Redhat Linux 7.2 / PC-AT 互換機

3.10.2. ツール

本システムの開発には下記のツール・ライブラリを使用する。

1. jdk-1.3.1_03
2. ant-1.4.1
3. log4j-1.1.4
4. wget
5. prescript
6. MySQL
7. Tomcat 4.0
8. Apache
9. perl
10. imdkcv

3.10.3. 開発の体制

各モジュールの開発体制は下記の通りである。

- ロボット

- 開発内容: 既存システムを流用
- 開発言語: perl
- 使用ツール: wget

- **書誌情報抽出エンジン**
 - 開発内容: 既存システムを修正
タイトル, 著者などを別個に抽出
 - Unicode 対応
 - 開発言語: perl
 - 使用ツール: prescript

- **参照情報作成エンジン**
 - 開発内容: 既存システムを修正
ロジックは流用
バックエンドに DB を利用する
論文及び参照の ID の発行方法を変更する
書誌情報・未解決参照情報の取得部分を新規に実装
 - 開発言語: perl or java
 - 使用ツール: wget (書誌情報の収集に利用する可能性有り)

- **マスタ書誌情報修正 UI**
 - 開発内容: 新規に作成 .
 - 開発言語: java
 - 使用ツール: tomcat

- **サーバ管理 UI**
 - 開発内容: 新規に作成
 - 開発言語: java
 - 使用ツール: tomcat

- **データベース**
 - 開発内容: フリーの RDB を利用
 - 開発言語: SQL, java
 - 使用ツール: MySQL

- **検索サーバ**
 - 開発内容: ほぼ新規作成
 - 開発言語: java
 - 使用ツール: tomcat

3.11. テスト

前節までで述べたシステムのテストを行った．テスト項目を以下に示す．

(1)統合システムに関するテスト

試験概要	試験手順	試験結果
データソース管理・新規作成・ローカル	ローカルデータソースの作成を選択し，登録を行う	入力したデータソースが登録される．Cron 設定が更新される．
データソース管理・新規作成・リモート	リモートデータソースの作成を選択し，登録を行う	入力したデータソースが登録される．Cron 設定が更新される．
データソース管理・設定変更・ローカル	ローカルデータソースの設定を変更し，登録を行う	変更したデータが登録される．Cron 設定が更新される．
データソース管理・設定変更・リモート	リモートデータソースの設定を変更し，登録を行う	変更したデータが登録される．Cron 設定が更新される．
データソース管理・本文データコピー	ローカルデータソースで本文データコピーを実行する	マスタディレクトリのファイルがデータディレクトリにコピーされる
データソース管理・本文データコピー	本文データのコピーの実行中に同じデータソースの本文をコピーする	コピーが実行できない
データソース管理・書誌情報更新・ローカル	ローカルデータソースで書誌情報の更新を実行する	書誌情報ファイルの更新処理が実行される
データソース管理・書誌情報更新・リモート	リモートデータソースで書誌情報の更新を実行する	データディレクトリ内の書誌情報ファイルが更新される
データソース管理・書誌情報更新・一括更新	全書誌情報一括更新を実行する	全ての書誌情報の更新処理が実行される
データソース管理・書誌情報更新・リモート一括更新	リモート書誌情報一括更新を実行する	リモートの書誌情報が全て更新される
データソース管理・書誌情報更新・ローカル	書誌情報の更新を実行している最中に同じ書誌情報を更新しようとする	更新が実行できない
データソース管理・書誌情報更新・リモート	書誌情報の更新を実行している最中に同じ書誌情報を更新しようとする	更新が実行できない
データソース管理・書誌情報更新・一括更新	書誌情報の更新を実行している最中に同じ書誌情報を更新しようとする	更新が実行できない
データソース管理・書誌情報更新・リモート一括更新	書誌情報の更新を実行している最中に同じ書誌情報を更新しようとする	更新が実行できない
参照データ管理・設定変更	参照データの設定を変更する	設定が変更される．Cron 設定が更新される．
参照データ管理・参照データ更新	参照データの更新を実行する	参照データの更新処理が実行される
参照データ管理・参照データ更新	参照データの更新の実行中に参照データ更新を実行する	更新が実行できない
本文データ手動追加・ローカルファイル	ローカルのファイルをアップロードする	アップロードしたファイルがデータソースディレクトリに追加される
本文データ手動追加・リモートファイル	URL を指定して論文を追加する	指定した URL の論文がコピーされる
書誌情報修正	指定した論文の書誌情報を修正する	データベースのデータが更新される
システム設定変更	システム設定の変更を行う	変更結果が登録される
システム管理者設定・新規作成	管理者を新規に登録する	該当する管理者が登録される
システム管理者設定・新規作成	既に登録されているアカウントに登録する	登録ができない
システム管理者設定・設定変更	アカウントの設定を変更する	変更内容が反映される
システム管理者設定・削除	システム管理者を削除する	指定した管理者が削除される
ログイン	正しいアカウントとパスワードを入力しログインする	管理メニューが表示される
ログイン	間違ったアカウントもしくはパスワードを入力しログインする	間違っている旨が表示され，ログインできない
ログアウト	ログアウトを実行する	ログイン画面に戻る

(2)検索インターフェースに関するテスト

項目

外部からのアクセスはできるか

トップページが表示されるか

トップページ

タイトルだけの入力で検索は出来るか

著者タイトルだけの入力で検索は出来るか

掲載誌だけの入力で検索は出来るか

著作年は過去から十分な未来まで対応できるか

データソースの選択は検索結果に反映されているか

表示順序の切替えは検索結果に反映されているか

検索条件を何も入力しない時, 入力を促すメッセージが表示されるか

1000 件以上対象があったら絞り込み検索を行うようメッセージが表示されるか

検索結果

総件数が表示されているか

表示件数以上の場合は複数ページに分けて表示がされるか

表示順序の切替えは検索結果に反映されているか

論文表示には、タイトル、著者、所属が正しく表示されているか

論文にチェックマークを付けられるか

検索結果のリストが複数ページに渡ったとき、それらを移動しても一度チェックした論文へのマークは消えないか

表示ボタンを押すとグラフの画面に移るか

1 つもチェックマークを付けずに表示ボタンを押した時、メッセージが表示されるか

論文タイトルをクリックすると、論文詳細情報画面へ移るか

再検索ができるか

再検索のフォームに前回の検索条件が残っているか

グラフ表示

縦の年代は表示されているか、また年代はあっているか

グラフ中にチェックされた論文は全て存在するか

グラフの下が古い年代、上が新しい年代になっているか

グラフの矢印の色は参照関係により異なっているか

グラフの矢印の方向はあっているか

矢印の色の意味は右下に表記されているか

ノードにポインタを置いたとき、zero-click は機能するか

初めてのユーザに zero-click を気づかせる仕組みはあるか

zero-click で popup する窓には必要な情報があるか

zero-click は galeon, netscape(gecko), IE, opera, safari(khtml engine) の各エンジンで表示できるか

H²のような遅い回線からアクセスしても zero-click の遅延は不便でないか

ノードをクリックしたときに論文情報のページが開かれるか

下フレームにチェックされた論文は全て存在するか

下フレームの論文タイトルをクリックしたときに論文情報のページが開かれるか

リンクをクリックすると参照関係詳細情報へ飛ぶか

論文詳細情報

各項目は埋まっているか

各項目が埋められないときは、何らかの処理がされているか

本文参照リンクはあるか

本文参照リンクは切れていないか

本文参照リンクが権利上切れている場合にはその旨を表示できるか

参照関係詳細情報

注目論文に対する参照論文、被参照論文が表示されるか

参照論文が表示されるか

被参照論文が表示されるか

引用箇所が表示されるか

表示順は 上から 参照論文、参照情報、被参照論文の順番になっているか

それぞれの論文の詳細情報へのリンクがあるか

3.12. 評価

前節で示したテスト項目について検証した結果，特に問題なく動作していることが確認された．付録にシステムの動作例を示す．

4. 実施計画書内容との相違点

- 書誌情報から表題，著者名等の抽出

これまでは，論文から抽出した書誌情報から，表題，著者名等の抽出は行っていなかったが，今回，書誌情報からこれらの情報の抽出を試みた．その理由は，複数の論文データから抽出された書誌情報の同定処理を行う際，もし表題や著者名が抽出されていれば，表記のゆれの大きいと思われる著者名や雑誌名等ではなく，ゆれのない表題部分を用いて同定処理を行うことで，同定処理の精度向上が図れると考えられるからである．また，ユーザが論文検索を行う際にも，あらかじめ要素毎に分けてある方がユーザの意図を検索に反映させやすく(検索に用いるキーワードが著者名を示しているのか，論文のトピック語なのか等)，それが検索精度の向上につながると考えられる．なお，書誌情報からは以下に示す項目を抽出する．また，抽出規則はMcCallumら[1]で用いられている隠れマルコフモデル(HMM)という機械学習の手法により獲得する．

著者名 : 奥村 学

第一著者名: 奥村 学

タイトル : 自然言語の意味的曖昧性の解消法

雑誌名 : 人工知能学会誌(ISSN09128085)

巻・号 : Vol.10, No.3

著作年 : 1995

ページ : 332-339

- SQL によるデータ管理

上記のように，扱うデータの粒度が細くなると，効率的にデータの管理や検索を行うためのより適切な方法を検討する必要がある．本プロジェクトでは，書誌情報や参照関係のデータ構造にリレーショナルデータモデルを導入し，リレーショナルデータベースのために標準化された言語 SQL を用いてデータ管理を行う．

- 抽出データの修正機能

書誌情報の抽出はシステムが自動的に行うが、抽出の失敗がある程度の割合見込まれる。しかし、状況によっては、このような誤りを人手で修正できる必要がある。例えば、データベース作成業者が論文データベースを作成する作業を支援するのに PRESRI を用いる場合が考えられる。この場合、PRESRI の抽出した結果を修正することで、ゼロからデータベースを作成するよりも少ないコストで、大量の高品質なデータ作成が可能になると思われる。また、PRESRI のデータを研究に利用する場合にもデータの修正の必要性があると考えられる。例えば、PRESRI のデータを引用分析研究に利用したり、インパクト・ファクタのような尺度で研究の重要度を測る場合、より正確な書誌情報に基づいてデータベースを構築しておく必要がある。以上のような目的への PRESRI の利用に対応できるよう、本プロジェクトでは、抽出された書誌情報が修正できる機能を PRESRI に加える。

- Unicode への対応

これまで作成したプログラムは日本語 EUC コードに対応しており、日本語と英語論文のみを処理対象にしていた。今後は、日英以外の言語で記述された論文も取り扱えるよう、プログラムを Unicode に対応したものに置き換える。

5. 開発成果の特徴

- 複数の論文データベースを統合的に利用可能な環境を提供

今日、多くの論文データが Web 上から入手可能である。我々は、これまでに Web 上に存在する Postscript および PDF 形式の日英論文データを収集して論文データベースを構築している。しかし、研究者が利用可能な論文データベースは、このような Web 上の論文データ以外にも数多く存在する。例えば、近年では、国際会議や学会の全国大会では予稿集の代わりに CD-ROM が配付されることが多い。また、研究者は、それぞれの所属する組織の図書館にある論文データベース等を利用することができる。この他に、学会や出版社の所有する論文データベースも利用できる。このようなローカルに存在する論文データを、PRESRI と統合的に利用できれば、非常に便利である。例えば、CD-ROM 中の論文と PRESRI 中の論文が参照関係にあれば、その参照関係をたどって、効率的に関連論文を集めることができる。以後、個人の所有する CD-ROM、図書館や出版社や学会が所有するデータベースをまとめて、ローカルな論文データベースと呼ぶ。本プロジェクトでは、このようなローカルな論文データベースと PRESRI とを統合的に利用できるようにシステムの構築を行っている。

● わかりやすい論文提示

ResearchIndex, Cora をはじめとする多くの引用文献データベースの論文検索インタフェースは、検索結果や参照関係にある論文をリスト形式で表示するのが一般的であったが、このような表示方法では、より大きな参照構造の中での個々の論文の位置付け(関係)がわかりにくいという問題点があった。本プロジェクトでは、論文間の参照関係をグラフで表示し、ユーザがグラフ上の論文アイコンにカーソルを重ねるとその論文の情報が、参照関係を示す矢印にカーソルを重ねると参照個所が提示できるようにしている。また、ユーザの用途に応じて、参照構造をズームインやズームアウト出来る機能も提供する。さらに、ユーザがグラフ上の論文アイコンにカーソルを停止させた時間と、その論文に対するユーザの関心の強さに相関があると考え、カーソルの停止時間に応じて、表示する論文情報の詳細度を変える機能を提供する。

6. 今後の課題・展望

- 本プロジェクトで開発したシステムをより多くの方に使っていただくために、普及活動を行う必要がある。まずは、以下の国際会議においてシステムのデモンストレーションを行う予定である。

41st Annual Meeting of the Association for Computational Linguistics
(July, 2003) 札幌

- 現在、東京工業大学学術国際情報センタで運用を検討しているサービスの中に、本システムを利用する方向で検討されている。

7. 開発分担

以下に開発分担を示す。なお、開発にはプロジェクト開発者、再委託先である DUO システムズの他に、東京工業大学総合理工学研究科博士課程の阿辺川武氏が加わっている。

- DB 統合システム設計：難波英嗣，奥村学，齋藤豪，阿辺川武，DUO システムズ
- ユーザインタフェース設計：難波英嗣，奥村学，齋藤豪，阿辺川武，DUO システムズ
- 書誌情報からのタイトル，著者名等の抽出：難波英嗣，阿辺川武
- ユーザインタフェース実装：難波英嗣，阿辺川武，DUO システムズ
- PRESRI の SQL への対応：難波英嗣，阿辺川武
- 結合試験：阿辺川武，齋藤豪，DUO システムズ
- PRESRI の広報活動：奥村学

8. 付録

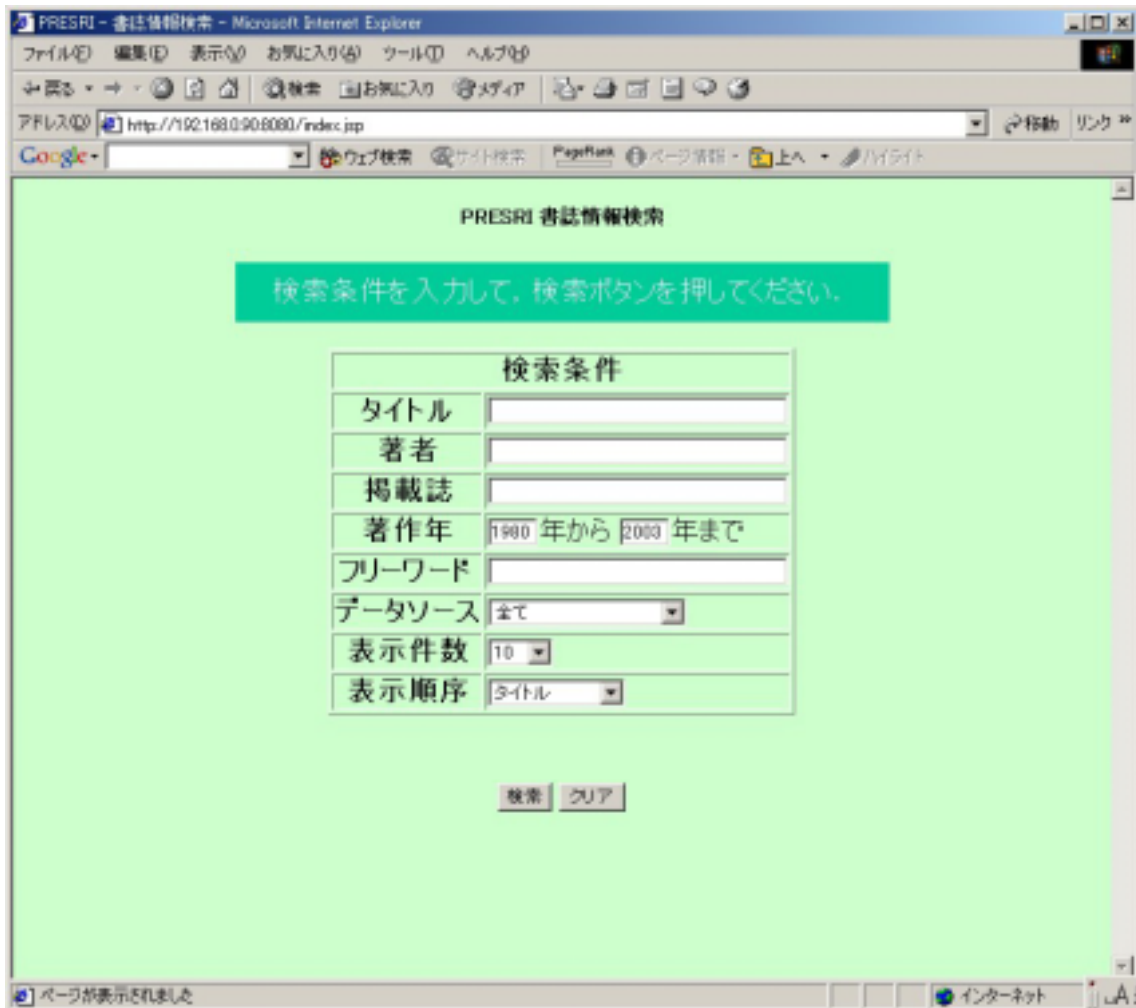


図 4 検索画面

これまでに開発してきたシステムと異なり、今回のシステムでは、タイトル、著者名、掲載誌にキーワードを分けて論文検索を行うことができる。

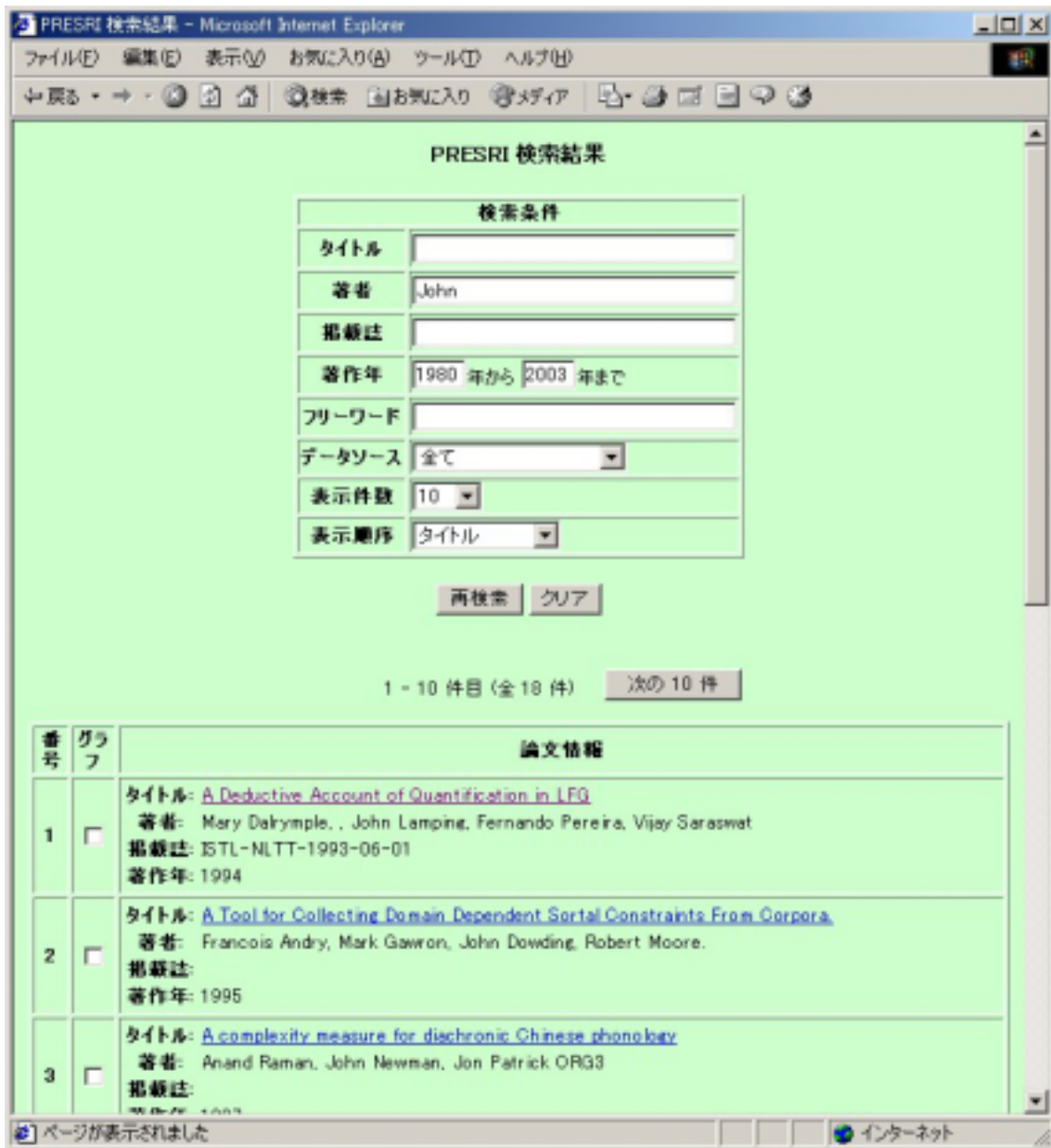


図 5 検索結果

図 5 は 1980 年～2003 年の間で著者名に John を含んだ論文を検索した結果を示している。図より，18 件の論文が検索されていることがわかる。

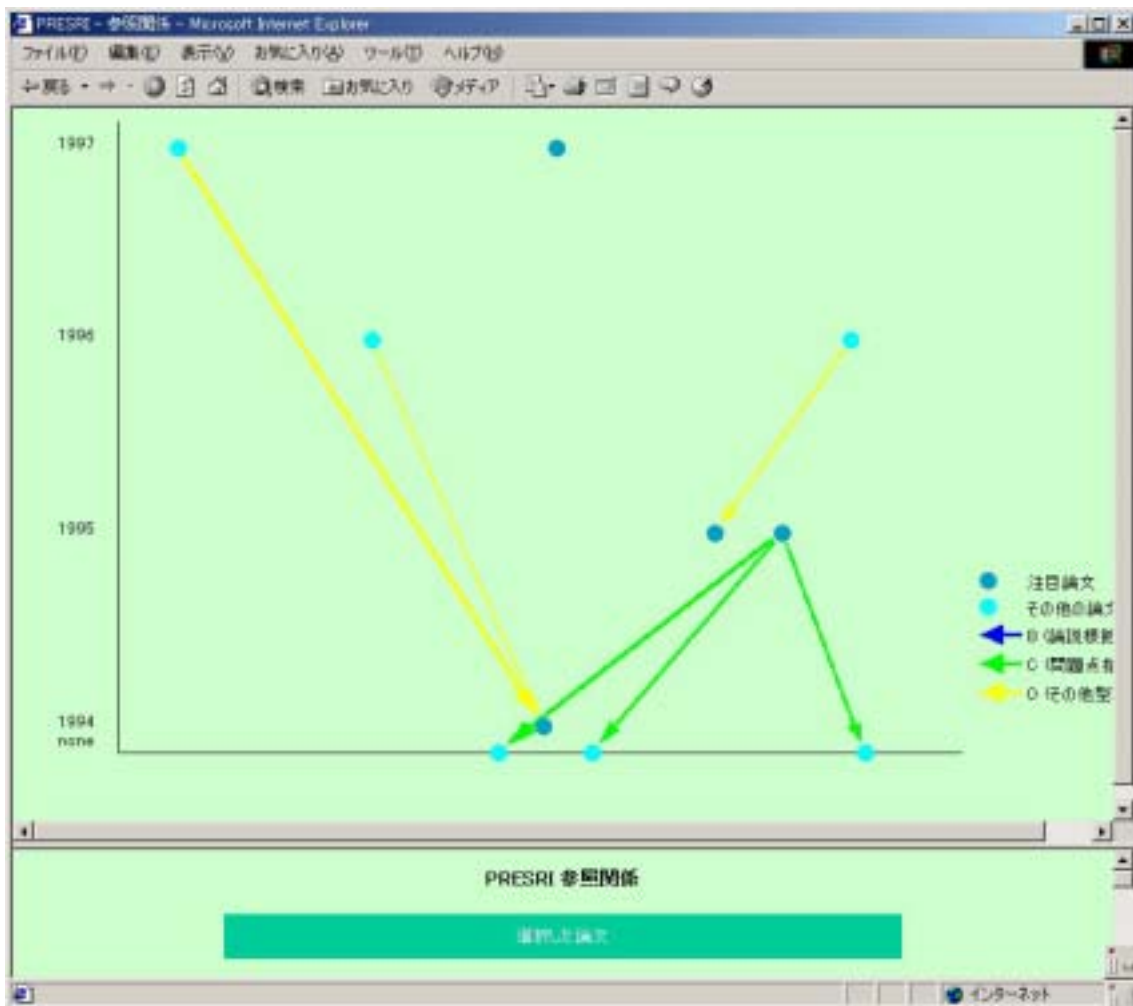


図 6 検索結果のグラフ表示

図 6 は、図 5 の検索結果画面において、チェックボックスでチェックされた論文と、それらに関連する論文を論文データベースから収集し、表示したものである。図において、「注目論文」は、図 5 でチェックされた論文を示している。「その他の論文」は、論文データベースから自動的に収集された、注目論文と関連のある論文を示している。また、図の矢印は論文間の参照関係を示しており、参照タイプ毎に色分けして表示されている。

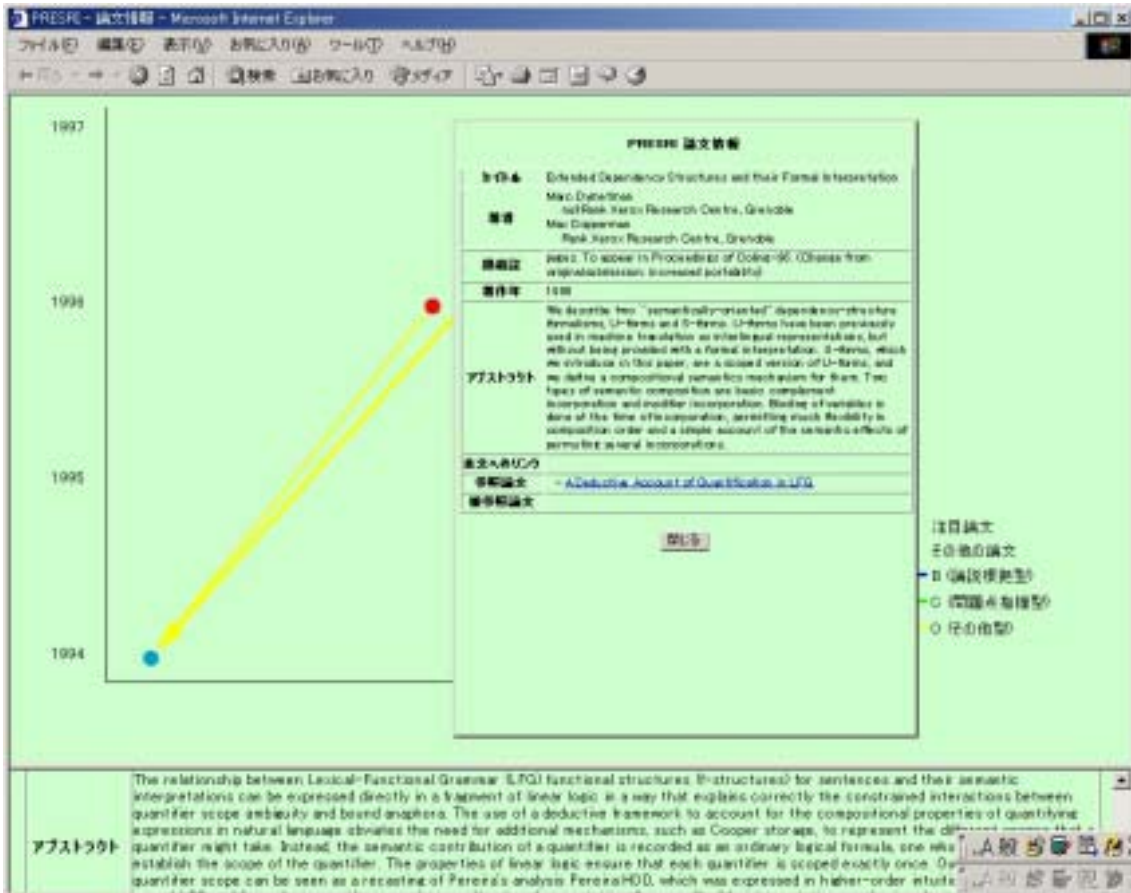


図7 グラフ表示における論文情報のポップアップ画面

図6中で個々の論文アイコン(ドット)にカーソルを重ねると,図7のようにその論文の書誌情報およびアブストラクトがポップアップ表示される。また,矢印にカーソルを重ねた状態でクリックすると,新しいウィンドウが立ち上がり,ポップアップ画面に表示されている内容がウィンドウ内で閲覧できる。

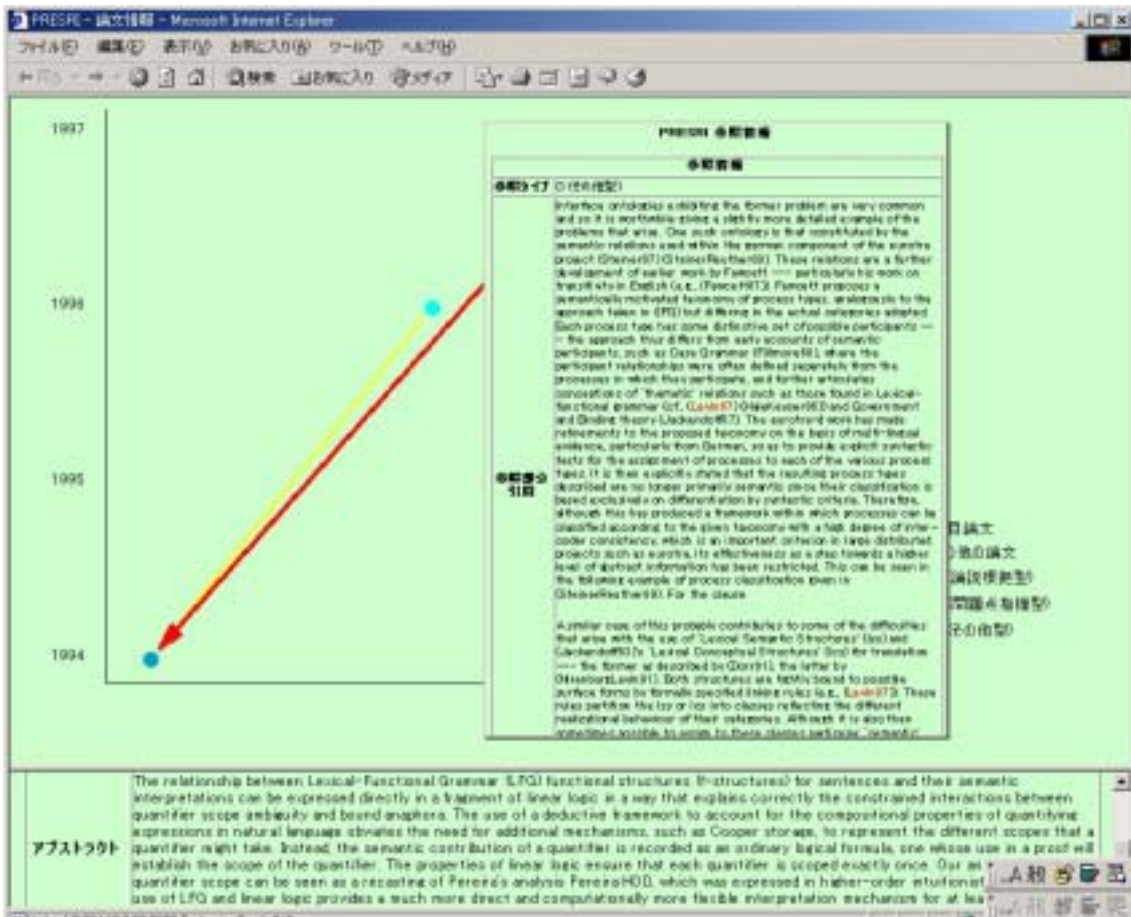


図 8 参照個所のポップアップ表示

図 6 中で矢印にカーソルを重ねると，図 8 のように参照個所がポップアップ表示される．論文中で複数回参照されている場合には，図に示すように，すべての参照個所が提示される．また，矢印にカーソルを重ねた状態でクリックすると，新しいウィンドウが立ち上がり，ポップアップ画面に表示されている内容がウィンドウ内で閲覧できる．

以下は、DB 統合システムに関する画面である。

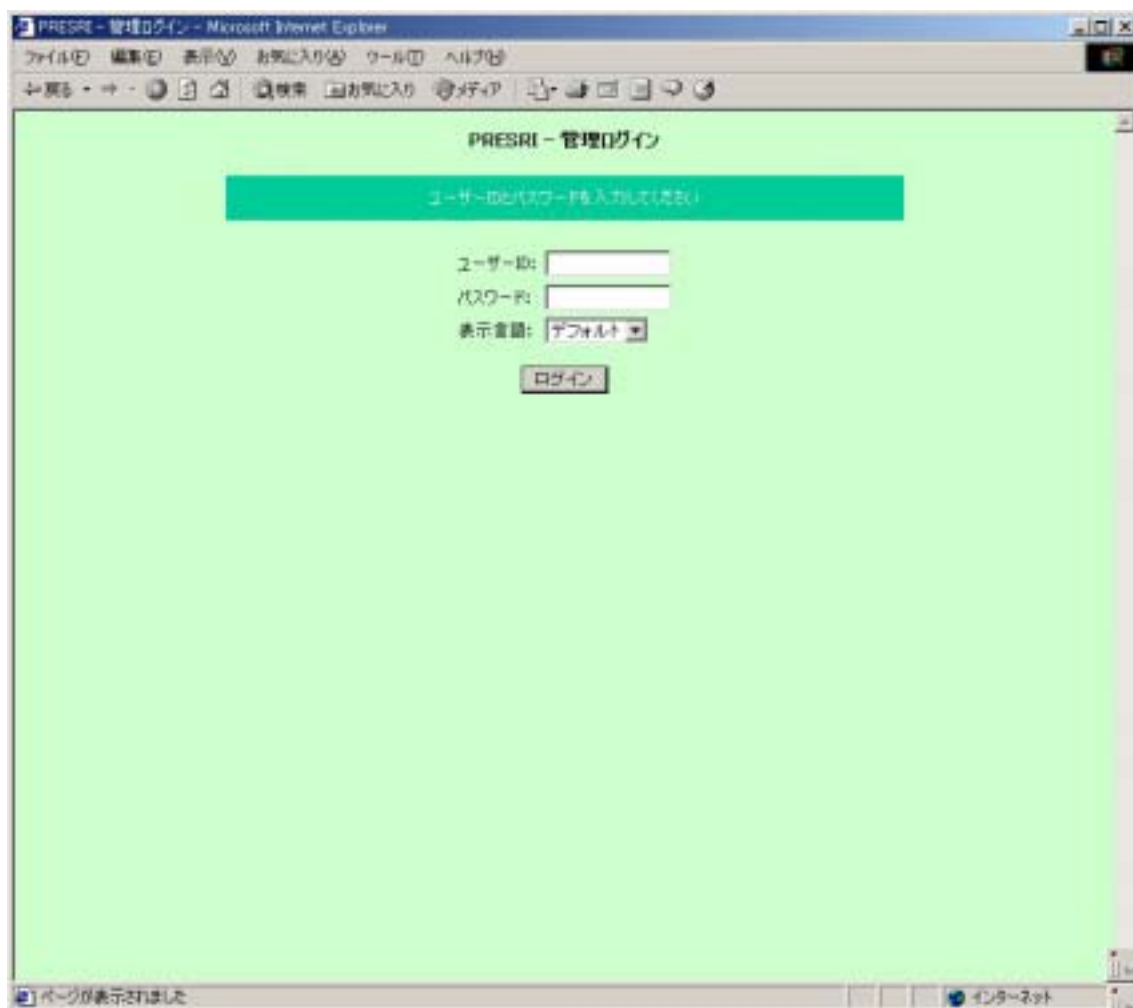


図 9 DB 統合インタフェース

図 9 は DB 統合システムのインタフェースである。複数 DB の統合は、Web ブラウザベースで行う。統合システムは、特定のユーザのみ利用可能になっている。ユーザは、ログイン時に、ID とパスワードを入力し、指定言語(現在は日本語と英語)を選択する。

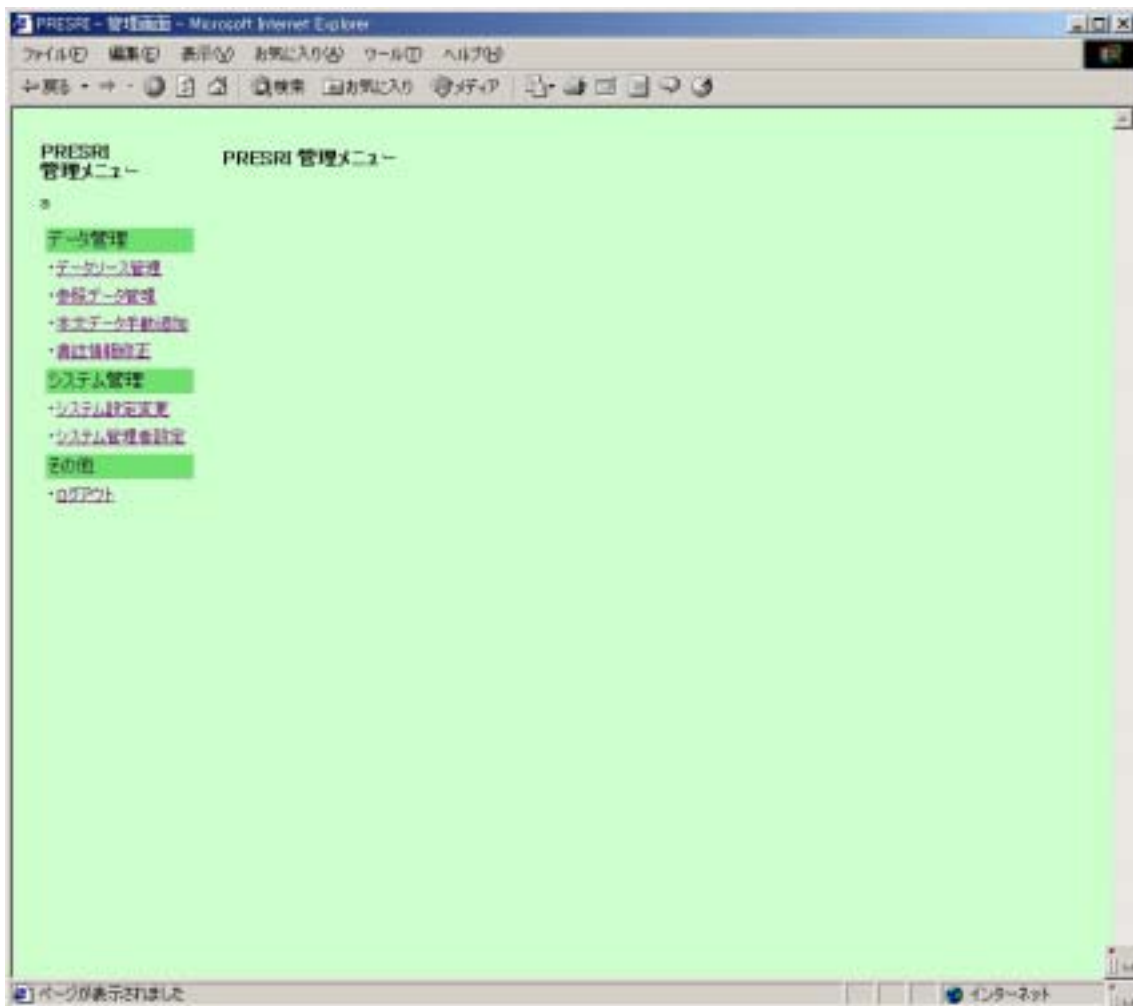


図 10 データ管理トップ画面

図 9 の画面において、ログインすると図 10 の画面が表示される。画面の左側に管理項目が表示されており、選択すると画面の右側に表示される。

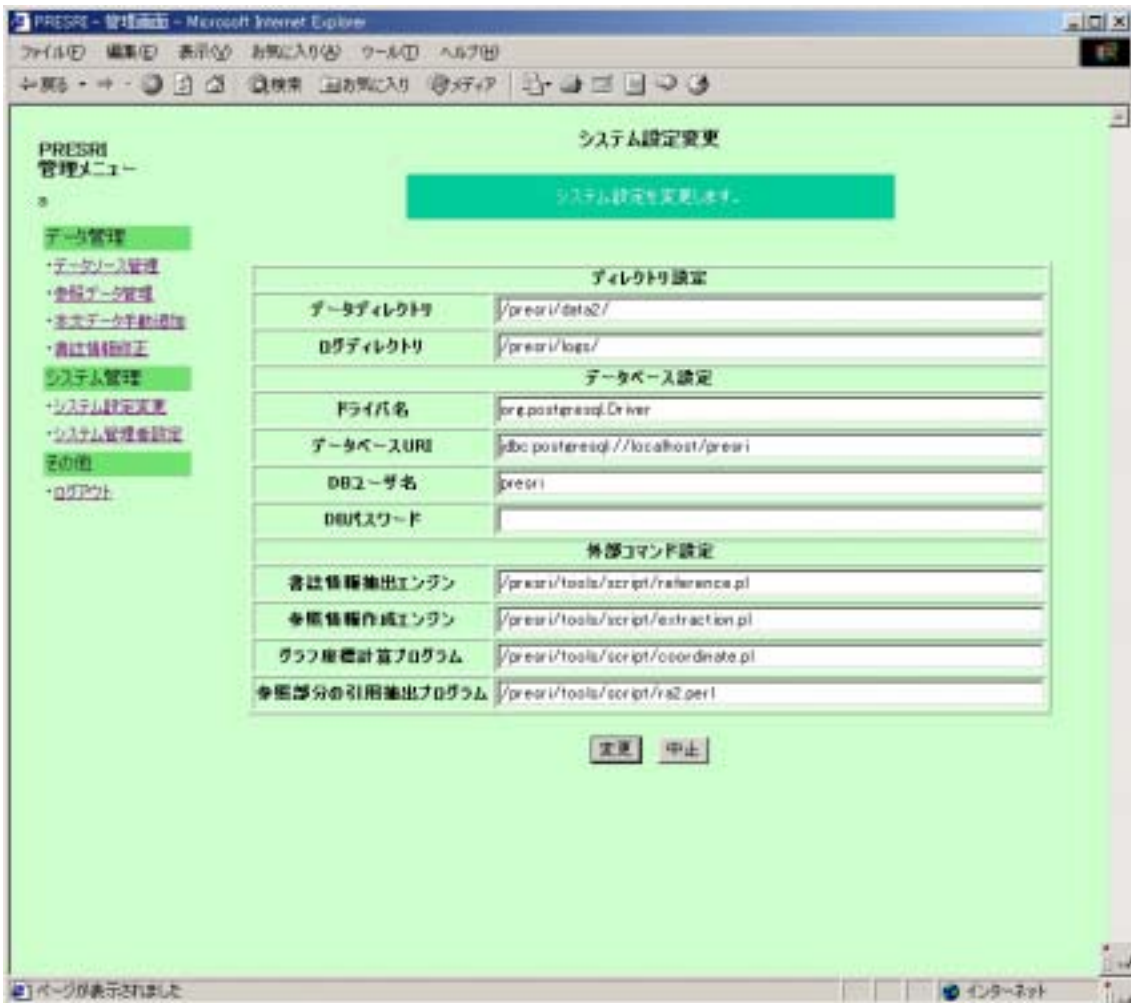


図 11 システム管理画面(システム設定変更)

システム管理 (システム設定変更) 画面において、PRESRI サーバ上に存在する論文データのディレクトリやログファイルの設定、SQL ドライバ、諸プログラムの所在の設定等を行う。

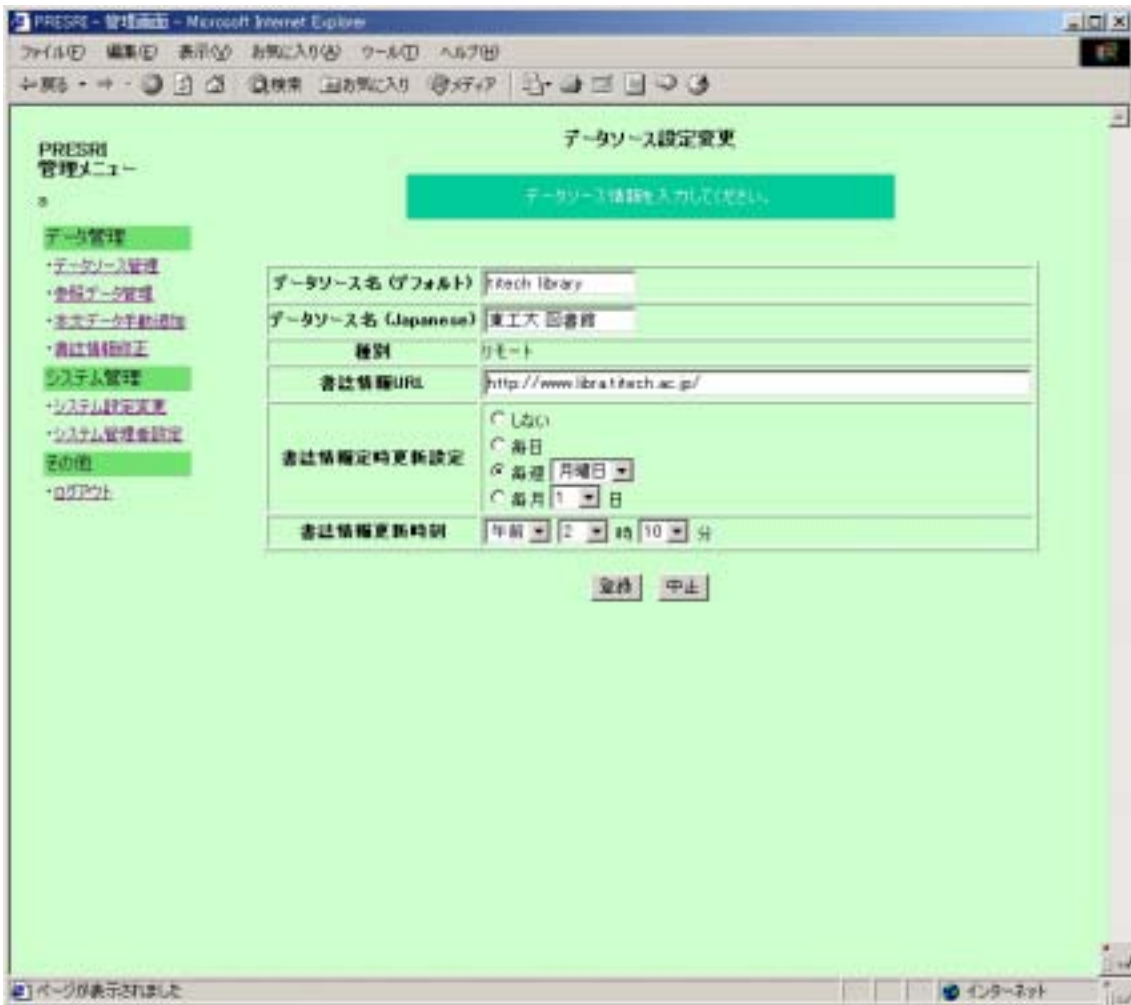


図 12 データベース(リモート)登録画面

図 12 はリモートに存在するデータベースの登録画面である。登録画面では、データベース名、データベースの所在(URL)、データベースの更新頻度を入力する。

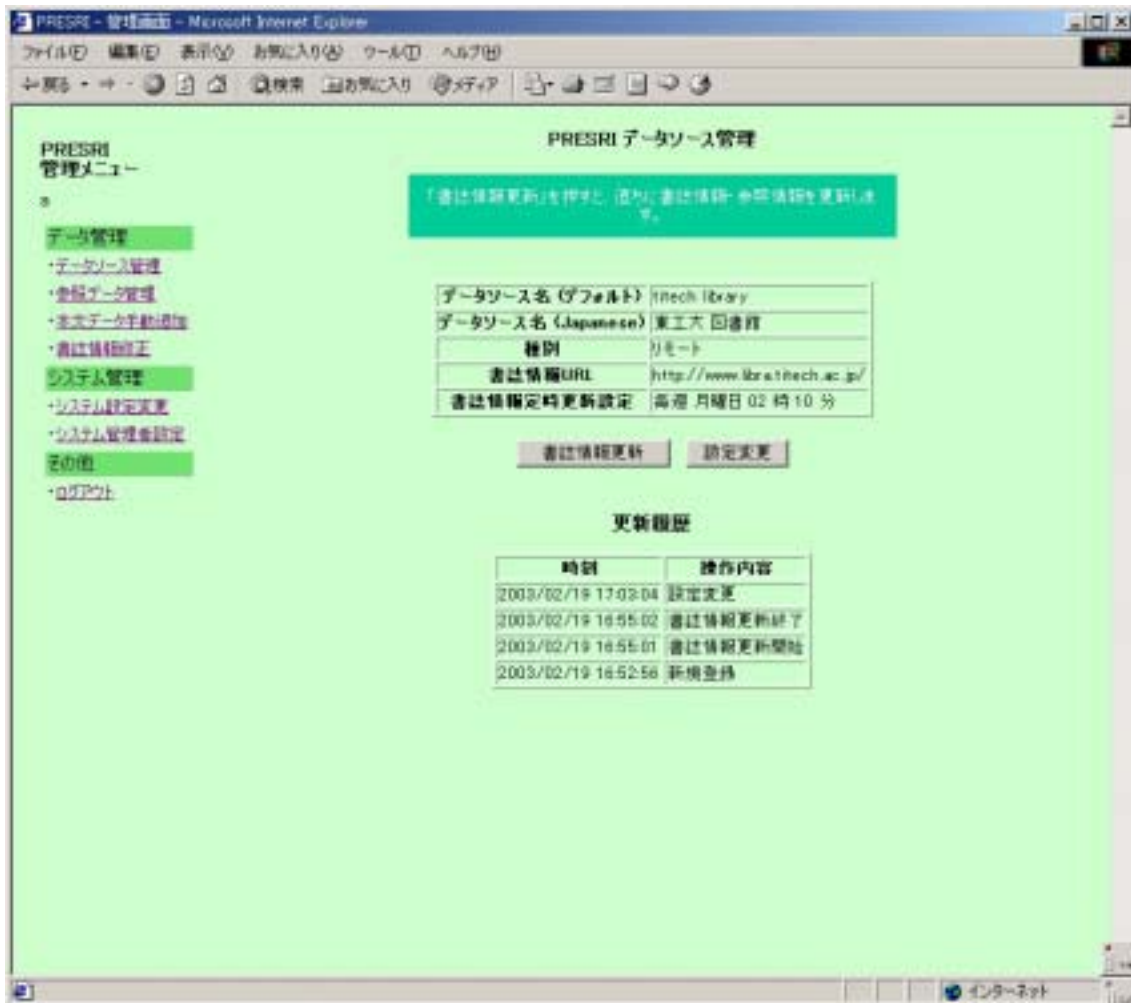


図 13 登録された DB の確認

図 13 は、登録されたデータベースおよびそのデータ更新履歴を示している。データの更新は図 12 の画面における設定で定期的に行われるが、図 13 の画面で「書誌情報更新」ボタンを押せば、手動でも更新を行うことができる。また、「設定変更」ボタンを押すことで、データベース名、更新頻度等を変更することができる。

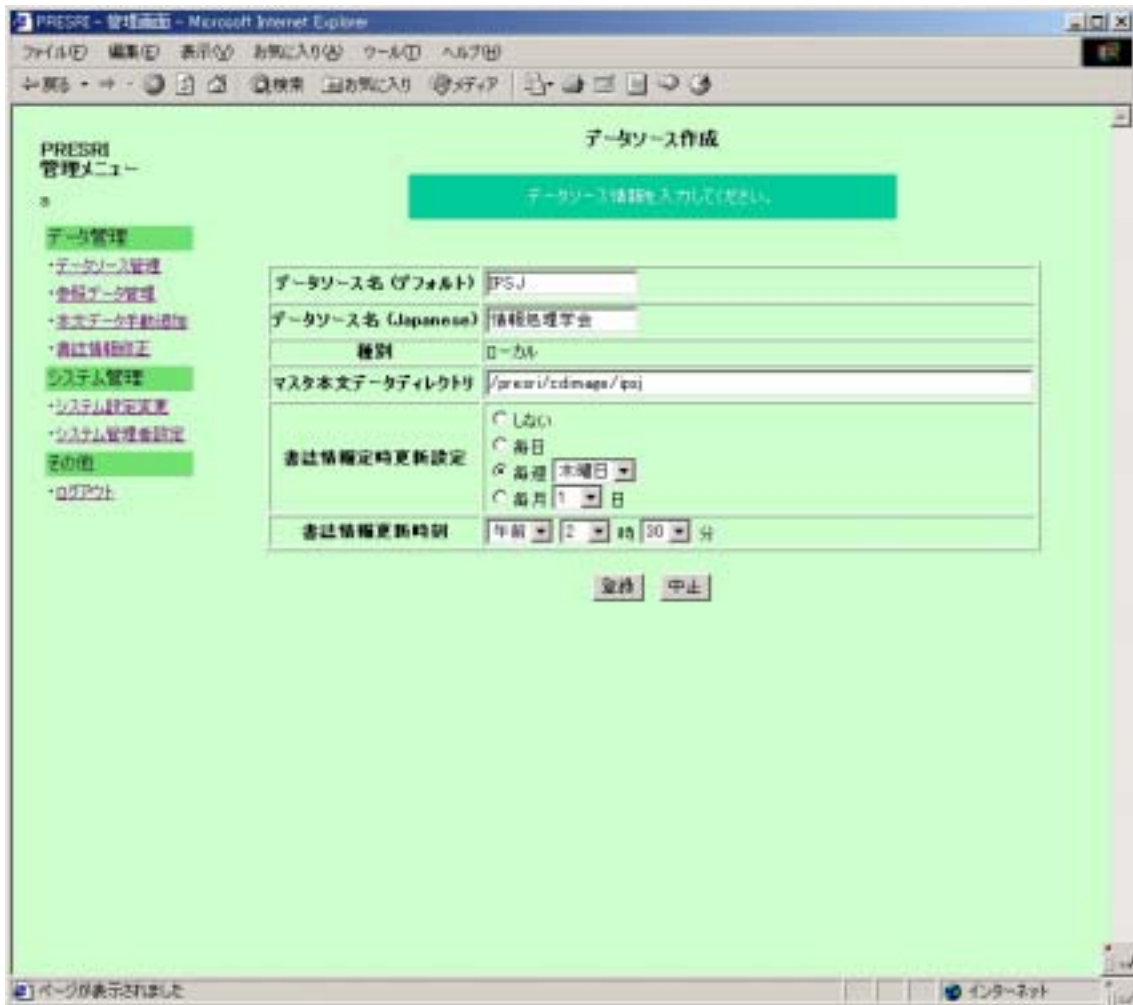


図 14 データベース(ローカル)の登録画面

9. 参考文献

[1] A. McCallum, K. Nigam, J. Rennie, and K. Seymore "A Machine Learning Approach to Building Domain-Specific Search Engines," The Sixteenth International Joint Conference on Artificial Intelligence, pages 662--667 1999.