

EADによる情報検索システムの構築

京都大学研究資源アーカイブ

五島 敏芳

京都大学研究資源アーカイブの五島です。12月から京都大学へ移りまして総合博物館に横付けされている所に所属をしていて、研究資源アーカイブを担当する事になっております、8年間ほどお世話になるかと思えます。そのおかげで実はこのプロジェクトとKEK 史料室さんにも今日の話をする内容ではご迷惑をおかけしてしまいまして、大変申し訳ありませんでした。国文研のデータベースの不具合にすぐ私が国文研に居た時は管理者権限持っていたのですが、今、その管理者権限を行使することが出来ずにいるので、後でこっそり行使してしまうんですけれども、始終不具合が起きてしまっているので大変申し訳ありません。ただそれでも尚且つ、こちらで報告させて頂けることに感謝申し上げます。

タイトルがプログラムの方で「EAD-情報検索システム構築」となっておりますけれども、ちょっと改めさせていただいて「EADによる情報検索システムの構築」とさせていただきます。同じ物です。

最初のこういう類の話というのは、今日纒々実は他の先生方がお話されてきておりますので、具体的な事を少し紹介して私の報告に変えたいと思いません。

EADという呪文の説明の前にアーカイブズの資料というものは、どうやって見つけてどうやって使える様にするかという情報検索システムの事について少しかだけお話しをしておきたいと思えます。

松岡先生が、図書館の世界でのお話をしたり、PORTA（国立国会図書館デジタルアーカイブポータル）のお話など画面で紹介されていたと思うのですが、そこで出てくるのは上澄みなんですね、Googleよりは性能がいいかもしれないけれど、全般的にある限られた領域のものを探し出してくるというものと五島は理解しております。有象無象を集めるといよりは見つけやすくする、その基盤というのは実はそれぞれのコミュニティで持っている情報がきちんと探し出せるようになっている事が大事です。

アーカイブズの世界の場合は見つけるためのツールというのが検索手段、finding aid と呼ばれるもので、日本では典型的には資料目録と呼ばれて来たものがそうだと思います。目録という言葉自体使われていますから、その説明はいいかと思います。

ダイブしなくても物が探せるようになるのが一番理想的ではあります。

なんですが、実は揚足を取るかのように、図書の資料目録を例に出すのは良くはないのですけれど、図書の資料と比べると、アーカイブズの資料というのは非常に扱いづらい、要するに面倒なんですね。

図書の資料の場合は、図書 1 冊だったらその 1 冊が何か 1 つの独立した情報を持っていて、その中の書誌事項は決まりきったものですが、アーカイブズの資料とはそうではないと言うのは今日もうこの場でも説明があったかと思います。出所があるとか、群として存在しているとか、量が多いとか、後は 1 個 1 個取り上げてしまうとそれぞれ断片的であったり、また同じ様な配布資料が何回も、それは会議としてよくないのかも議事録の作り方が悪いという事なのでしょうけれど配布資料で同じ物が出てくる、重複する部分も現実には沢山あると思います。そうなると 1 個 1 個の中身だけを見ていても仕方がなくて、前後の脈略であるとか構造も把握しないといけない、そうしないとアーカイブズの資料は理解出来ないというのが前提にあるとするとその為の探すツールというのは図書の目録では足りなくて、前後の脈略の説明が中に組み込まれていないといけないと言う事があると思います。

階層に関して、具体的にどの様な階層があるか説明を簡単にしていきたいと思います。

アーカイブの資料は出所があって、その単位で捉える一番上の collection、fonds と呼ばれるまとまりなんです、その中の例えば組織であれば分業の体系、人間であれば機能等で分けられていると思うんですね。

まとまりとしての理解の仕方は一番上の 1 個だけ理解しても勿論良いと思うのです、まとまり全体、段ボールで何箱かあって、中を開けていないがどこから送られて来たというのは判る状態、それが 1 つ。

中を開けた場合、括りがあり、括りのまとまり、或いは意味で分類されていてここからここまでが何か判っている物がそれぞれあったりする。そういった物を一々記述しないといけないという面倒な事があります。今言った

パターンだと大雑把に分けられている、箱が何箱かあって箱の単位で分けられているでもいいですし、出所があって sub fonds、sub group と言うのは1つのパターンですし、それ以外にももう少し中が判ったというのがあったとすれば、そこまで区画として分けられた様な何か捉え方も出来るでしょう。それ以外に片端から1点ずつリストにしていった等の捉え方もあるかもしれません。全ての階層でという事があるかもしれません。そういった幾つかの捉え方もあるし、図書の場合は1冊に1レコードと対応するのですがアーカイブの資料はそういったデータの捉え方が出来なくなってしまう、collectionがあるとすると、例えばこれは五島所蔵史料なのですが古本屋で買ってしまった物1個これがcollectionだとしたら、これに対する概要のデータ1個記述を作ると1個記述レコードが出来ますね、これで終わりにしてももちろん良いし、中を1枚1枚取り出してレコードを作ってもいいと思いますが概要データでレコードは確実にある訳です。その次に札が見えるのでしょうか1とか2とか3とかありますが、纏りによって大よそ意味的にも区画的にも分けられていると考えて、ある series だと仮に捉えたとすると、series の概要データが3つ出来るわけですね、一応階層構造この様に捉えること出来るのですが、足すと1個の資料のデータに対して合計すると4つ出来ます、目録レコードという言い方が正確かわかりませんがレコードが4つ出来る、作業の進行によってデータが増えるという不思議な状態が起こるその集積体がアーカイブの資料の目録なのです。

やっと EAD の話なのですが、EAD とは何かというと図書の目録の場合は、図書の目録を作るにあたり、目録の規則が前提にあります、その規則に則りデータを書いていく、その書いていって図書だったら OPAC(Online Public Access Catalog) に載せます。その段階で加工しないといけない、加工する上に企画も揃えないといけない OPAC の場合は ISBD(International Standard Bibliographic Description) 記述があり国際標準があります、国際標準は考え方なので規則までは細かく定められてはいない、それについて書誌的事項としてどうとれば良いか目録規則として存在しています。日本目録規則でもいいですし NCR でもいいですしあるとすると、それに対して OPAC で具体的に実現するのは MARC(The Machine-Readable Cataloging format) というのがあります。

(資料7) 表下側に ISBD に対して MARC とありますが、そのアーカイブズ版みたいなものです EAD というのは。

アーカイブズ版の機械可読目録、電子化するための構造の標準、デファクト標準と言うのが EAD(Encoded Archival Description) と言うことが出来ます。ほかに EAC や EAG など存在しますが細かく説明しだすと時間が足りなくなってしまうんですけど、図書の目録で著者であるとか編集者から本を探したいという場合もあると思います、その時に著者の情報が別途入っていて、ペンネームと実名とバラバラになっているとしてそれが同じ人という情報があると、ペンネームの方で引いても実名で書いた本や論文が出てくるようになるので便利ですね、そういった人に対する情報をまとめているのが実は典拠レコードといいまして、そのアーカイブズ版、図書館典拠レコードの規則は忘れてしまいましたが、MARC については書誌の情報以外に典拠レコードの何か部分があったかと思います。それに対応する物が Encoded Archival Context (EAC) と言うことが出来ると思います。Encoded Archival Guide (EAG) というのは何処にどんな史料があるか、どういった所が持っているか、持っている所に関する大雑把な情報をまとめたもので、つい最近 α 版の DTD の中身を解析することが出来たのですが、今は説明がまだ出来ないので申し訳ありません、これは省略します。それは何で書かれているかと言うと XML、または元々は XML が基になっている標準汎用マークアップ言語 (SGML) で書かれる物です。EAD というのはそのセット、部分集合サブセットという言い方も出来るかと思います。DTD 文書形定義またはスキーマ (schema) 文法や記号等に規程、制御している、XML / SGML を使わなくても EXCEL 等を使う方が便利じゃないと言われるのですが、先程の階層的な把握と言うのをやらないといけないので、その為にはデータを入れ子に扱わないといけないのですね、具体的に例を挙げないとピンと来ないかもしれないのですが、どう階層制を填補するのか実は EAD の情報検索システム構築で一番大変だった所でもあります。それを実現できる枠組みと言うことだけ、大雑把に頭に入れておいて頂ければと思います。

この間にタグのセットがあるのですが、EAD というのは 3 つ位の大雑把な情報のまとまりがあって、その中にさらに細かい情報がセット出来る様になっています。ですから単純に CSV の様に、或いは EXCEL の縦横の表の様

にどこの列に何を入れ行が増えていくイメージとは違うのですね。

それもやろうと思えば出来なくはないのですが、どこかに階層か何か、何番目に何処の中に入っていると情報を別途入れてやらないと CSV の情報などに開く事が出来ない、と言う場合は先程階層的に 3 つか 4 つ位の把握の仕方がある事を申し上げましたけれど collection のレベルと series のレベルだけの情報を引っ張り出す記述をする事も出来ますし、collection のレベルと一覧のリストの組合せる事も出来ます、フルのセットというのもデータとして収める事は出来る、先程階層的な把握を言ったその階層的な幾つかのパターンを全てカバーできる枠組みが実は EAD の枠組みです。

と考えると何故 EAD なのかアーカイブズの資料の目録の難しさを克服出来たものだからと、特にコンピューターの資料目録、オンラインの資料目録で実現出来るから EAD を採用するのが良いのではないかと言う事です。

実はこれだけではアーカイブズの資料足りないんですね、メタデータとして使うには EAD だけでは不足しますが差当たって、何処にどんな情報があるか EAD に関して言えば、目録の情報に関して言えば EAD があれば足りるだろうと言う事です。

SGML、XML であることは重要な理由になりますし、何より CSP などだと、EXCEL の表 1 セルの中に本で言うと何十頁となる記述を 1 セルに収めることが出来るか、収められないですよ、今は出来る様になってしまったのですが、前は 256 文字しか確か入らなかったと思います。

実は図書館の資料目録を電子化した時 OPAC は、その制約にだいぶ縛られていました。それを壊して超えることが出来たのは実は EAD というより XML だった。尚且つ、入れ子の表現が出来るという事は、電子的検索手段としてマルチレベル記述というのが可能になる、つまり先程の階層的にデータを取扱うことが出来るようになりました。またそれだけではなく、これは EAD というより本文を電子化するための TEI (Text Encoding Initiative) という人文系で使われている物が別にあるのですけれども、そこの中で電子目録、電子書類として取扱えるような部分があり、先程のですと EAD ヘッダーと呼ばれる資料目録のバージョン管理のような事がそこで出来るのですね、というものを兼ね備えているので EAD を使うのがよろしいのではないかという事です。

では実際に OPAC みたいに出来るかという点と相当困難が最初ありました。

コンピューターがそこまで使えない、私が使えないというだけではなく、コンピューターにも駄目な部分もどうしてもあるので、コンピューターにさせる作業を限定的に考える事をやらないといけなかったと言うのは1つあります。

もう1つは、情報検索システムはアーカイブの資料で重要な脈略であるとか構造を、ズタズタにするのが大前提にあるのです。検索結果として現れてくるものは、どれか1個出てくるわけですね、それが1つの単位になっていて、脈略や構造の提示などとは完全に相反するものだと、そうなると困るというので、階層的な提示というのが出来るか確認をして次にキーワード検索とその表現へという形で話を進めて行き、この場でも使ってもらえるような所に、やっと至ったという話です。

その組合せが、参考文献で私が「記録と史料・第18号」に記事として書かしてもらった物の中にもありますけど2つのコンポーネントがありまして、表現の所、それと後は実際の検索システムという形で組んであります。実は EAD は資料目録のデータを作ってしまうと、一種の画一化・企画化なので、機関とか、機関の中の資料群の違いは、全部吸収して横断的に検索する事は、自動的に可能になるのですね。

(資料13) 全体像をこの様に出しましたが、上側に外部サーバーとありますが、ある EAD の検索システムが1つ動いていたとして、それとは別の所に EAD のデータ上のほうに置いてあったとしても、集めてきて中に取り込んで別の資料だと、別の資料 collection ですよという形で自分の所で検索システムがあれば収蔵資料目録のデータを併せて検索する事でそのまま総合目録・ユニオンカタログと同じ事が実現出来てしまいます、それをここで研究機関、研究資料のアーカイブズについても行うという事で、国文研に仕組みを入れてもらっていたので、そこを使う予定だったのですけれども急に居なくなると本当にご迷惑をおかけしました。

どの様な内容かを紹介し、全文の表紙は要するに概要があって、中身の一覧等があるのですが全部プリントアウトして行くと冊子形の資料目録と同じ物が出来上がります。それを Web のブラウザの画面上で表現するのが全文の表示の所です。これはあった方がおそらく良いでしょう。前後スク

ロールして見る事が出来ますから。

もう1つ、(資料21) そっけない画面が出てきてしまいましたが、本当はもっとインデックスを充実させないといけないと思っております、アルファベット順であるとか、事項の主題等でそこをクリックすると検索結果がズラッと出てくるなど、一覧が本当は必要であると思っております。以前そういったご指摘を総研大アーカイブズプロジェクトからも頂いた事があります。参加している基盤機関だけの collection 一覧を出してほしいとかそういったのがありました。出来るようにはなっているのですが今それをサーバーに登録する事が出来ないのでも止まってしまっています。

詳細検索であるとか簡易検索の項目などはこういったものかなというのを挙げておきましたので今は省いてしまいます。

簡易検索から(資料23) 仮に「物理」というキーワードを入れて検索をクリックすると、こういった検索結果(資料24-26) が出てきます。国文研の史料情報共有化データベース、全国アーカイブズ総合目録の検索結果の所では約800の機関が参加して頂いているのでこのように出てきますこれをスクロールしていくと、中にこの様な画面が出てきます(資料26)。

赤の矢印が出ている所が実は検索結果として当たった所なのですが、例えば早川幸男先生寄贈史料で、箱B303 bの中のある series というより具体的には file が出てきました、さらにその file の中に item をとっていて、その item の中でこの1点が該当しましたとこのように階層的に表示出来るようになりました。collection レベルのデータも series レベル或いは sub series レベル、file レベル、item レベルと順々にこの様に出せるようになっているのが、情報検索システムの中で階層的な表示を両立させる苦労した点でした。どれか1つ選ぶと詳細な表示が出てきます。(資料27-28) ここも少し工夫をした所で、上にどんな物が入っているか分かるようになってきました。また画面左側にナビゲーションのパネルを作成しまして、例えば今概要の所へラインを引きましたが、概要を選択すればその上位の属している collection レベルの情報に移ることが出来ます。

それ以外別途に画像・音声・映像等がこの資料については有ると、予めまとめておく事も出来ますし、それを選ぶと、ここでは別のメタデータ、データの為のデータを使わないといけないのですが、EAD以外のMETS (Metadata

Encoding & Transmission Standard) というものを使っています。画像のファイルとなると、ある冊子に頁があるわけですが、その1頁1頁が1つの本だとまとめておくデータも別途必要になるので、それもXMLでコントロールすることをやってきました。ここまでは現在総研大のアーカイブズプロジェクトあるいはKEK史料室では実現出来ていませんが、やろうと思えば出来ます。

例えば古文書の場合、タイトルが表示されている所をクリックすると説明の所に戻るというように、行き来出来るようにしてあります。

機関を超えて実現出来れば自然と全国アーカイブズ総合目録、アーカイブズユニオンカタログが出来上がるだろうと考えております。今の所国文研でしか実現できていないのですが、仕組みとしては、私の所属している京都大学研究資源アーカイブで同じ物が出来るかと思えます。差当り学内版ですが。

全国アーカイブズ総合目録の概要など言葉で示している物は今は省略したいと思えます。

もしここにいらっしゃる方でまだ参加されていない所がございましたらご紹介いただければと思います。具体的に所蔵資料を、公文書でなくても研究資料のまとまりなど色々なケースがあるかと思えます1つの枠にまとめて載せてみると先程のようにひっかける事ができます。1個2個であっても。しかしGoogleのように有象無象がでてくるノイズの多い検索結果ほど大きいことはないと思えます。実務的に行っていくと実感していただけると思えます。そこで面白い実例をみつけたら上の人に訴えかける事が出来るかと思えます、そのためのツールとして使ってみる、その為にはまずここに入力するか、同じ規格の物を作らないといけないので是非ご参加いただければと思います。

最後に宣伝ですが京都大学研究資源アーカイブにて学内で同様の事を行っています。ここでの横断検索を考えていまして、出来れば総研大のアーカイブズプロジェクトで出来上がっているようなEADのデータを横断検索対象として繋げさせてほしい、実証実験も予定されておりますので、国立公文書館さんにもお願いをしなければなりません、ご協力いただければと思います。学内版で2009年4月始動予定です。体を壊さない程度に両方とも頑張りたいと思えます、以上です。

EADによる情報検索システムの構築 アーカイブズ情報共有の実験

五島敬芳
京都大学研究資源アーカイブ
Archivist for Digital Collections
京都大学総合博物館
h.gotoh@inet.museum.kyoto-u.ac.jp

『記録管理とアーカイブズ』(3)

はじめに

- アーカイブズの資料は、見つからなければ、保存も利用もできない。
- その、資料を見つけるための工具(ツール)が、「検索手段」finding aidsである。
日本では、その典型が(伝統的に)資料目録といえる。



イラストは、オーストラリア国立文書館ブックレット“Keep the knowledge - MAKE A RECORD”より。

『記録管理とアーカイブズ』(3)

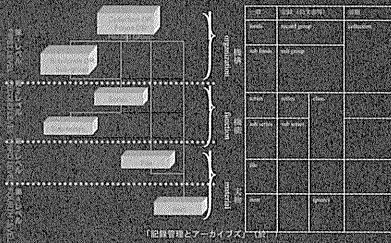
アーカイブズの資料目録の難しさ

- アーカイブズと他の資料とのちがいは、
 - (たとえば図書資料と比べると)
 - 出所の存在、再として存在、欠落、階層的構成、同一複製物なし、断片的な内容、重複冗長...
- 資料の(内容)だけでなく、(脈絡) (構造) の把握が重要。
- (叙述的) 記述中心、その提供順序は(大から小)
 - cf. 図書の資料目録は、平板なリスト(一覧表)だけでも足りる。

『記録管理とアーカイブズ』(3)

注) つぎの文脈の範囲に参照、OAC Working Group, OAC Best Practice Guidelines for Encoded Archival Description Version 2.0, CDL, 2005

資料群の階層構造モデル



『記録管理とアーカイブズ』(3)

資料群の階層的把握



- 群から部分へ
- 時間とともに記述データは増加

『記録管理とアーカイブズ』(3)

EADとはなにか?

- アーカイブズの検索手段(資料目録)を電子的身符化のためのデータ(事実上の)国際標準(規格)。
- EAD, Encoded Archival Description (符号化永久保存記録記述)
- EAC, Encoded Archival Context (符号化永久保存記録脈絡)
- EAG, Encoded Archival Guide (符号化永久保存記録(収集者)便覧)
- 記述の間隔標準頭との対応:
EADは、アーカイブズ版MARC (機械可読目録)

文書館	記述	ISAD	EAD
	典拠レコード	ISAAR	EAC
	収集機関	ICA-ISDIAH	EAG
	機産	ICA-ISDF	?
図書館		ISBD	MARC

『記録管理とアーカイブズ』(3)

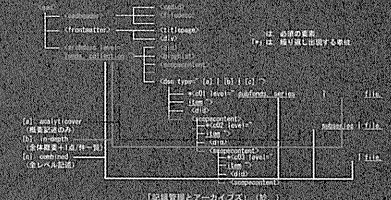
EADとはなにか? (contd.)

- 標準マークアップ言語SGML/XMLのサブセット。
- DTD/Schemaによる文法・記法を規定・制御。
- 文書型定義DTDは、1998年8月第1版公表、2002年10月2002版公表。(要素・属性の出現構成を提示。)
- スキーマ(Schema)は、2006年12月公表(2008年4月改訂、W3C Schema原とRelaxNG受)。 (使用可能な属性値をも提示。)
- 「収録目録登録簿として知られるアーカイブズの検索手段を符号化する」。
- 現在2002版(1度改訂)。

『記録管理とアーカイブズ』(3)

EADデータの構成と類型

- 記述水準・深度と各記述部分の組合せ。



『記録管理とアーカイブズ』(3)

なぜ EAD なのか？

- アーカイブズの資料目録の難しさを、克服できた。
 - コンピュータでの資料目録（電子的検索手段）で。
 - SGML/XML であること。
 - 大量の記述データを扱えること。
 - 電子的検索手段としてマルチレベル記述を表現できること。
 - データの階層的構造化が可能。
 - 入れ子状のデータ配列の実現。
 - 電子的検索手段としてだけでなく電子記録・電子書写として扱うための枠組みがあること。

【記録管理とアーカイブズ】 (続)

アーカイブズの情報検索システムへの挑戦

- EAD にいたるまでの困難。
 - コンピュータ利用を、あえて限定的にしか行なうことも。
- 〈検索システム〉の特性による矛盾。
 - もともと検索システムの結果表示は、縦書きを寸断してデータを断片化する。（←縦書き構造の表示と相反する。）
- 戦術；まず基本的表現（階層的提示）を確保し、つぎにキーワード検索とその表現へ。

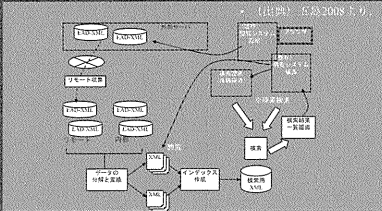
【記録管理とアーカイブズ】 (続)

EAD/XML 検索システムの紹介

- 資料目録 EAD/XML 化の実践例
- (A) 基本的表現：配信システム
 - 1) 資料目録全体の全文表示。
 - 2) 資料目録各部分の階層的表示。
- (B) 情報検索システム
 - 1) 特定の資料群の検索。
 - 2) 複数の資料群の検索。
 - 3) 複数の収蔵者の資料群の検索。

【記録管理とアーカイブズ】 (続)

EAD/XML 検索システムの全体像



【記録管理とアーカイブズ】 (続)

EAD/XML 検索システムの紹介

- 資料目録 EAD/XML 化の実践例
- (A) 基本的表現：配信システム
 - 1) 資料目録全体の全文表示。
 - 2) 資料目録各部分の階層的表示。
- (B) 情報検索システム
 - 1) 特定の資料群の検索。
 - 2) 複数の資料群の検索。
 - 3) 複数の収蔵者の資料群の検索。

【記録管理とアーカイブズ】 (続)

(A) 1) 全文表示

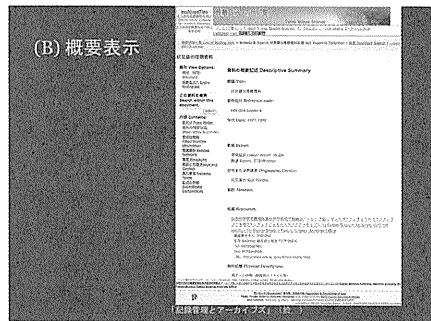
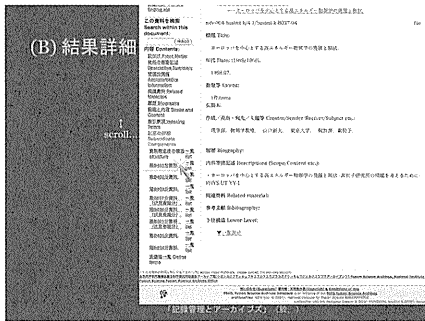
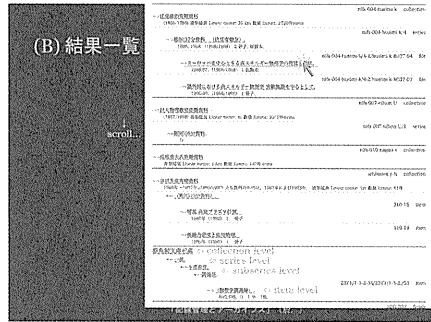
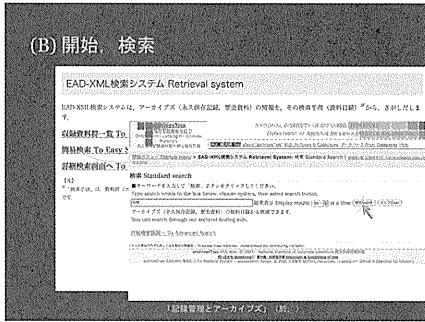
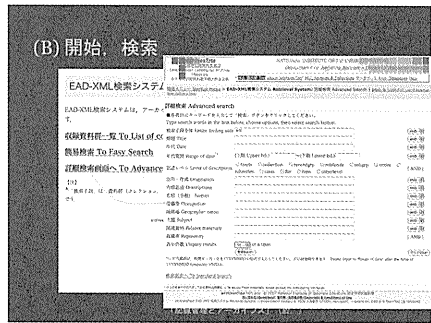
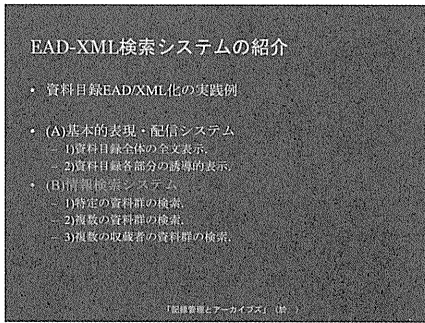
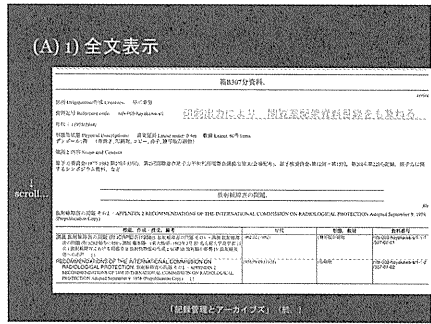
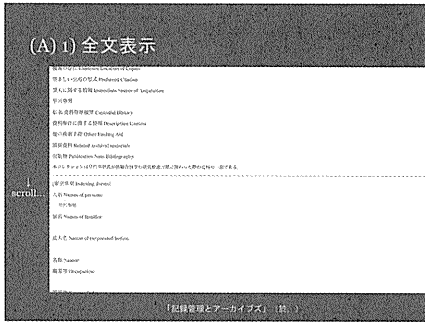
【記録管理とアーカイブズ】 (続)

(A) 1) 全文表示

【記録管理とアーカイブズ】 (続)

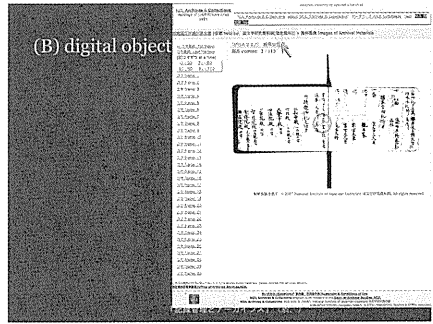
(A) 1) 全文表示

【記録管理とアーカイブズ】 (続)

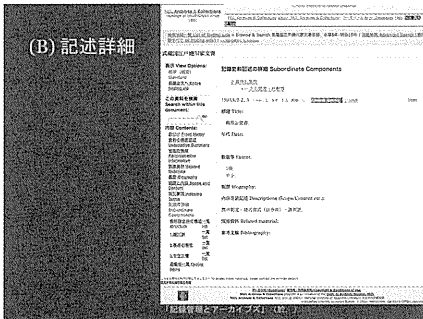




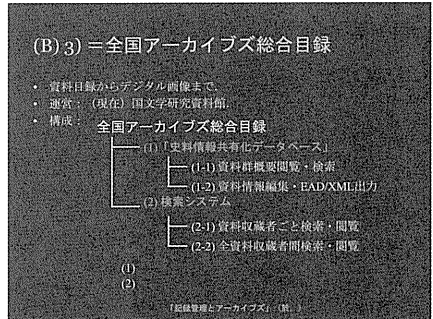
(B) 部分一覽



(B) digital object



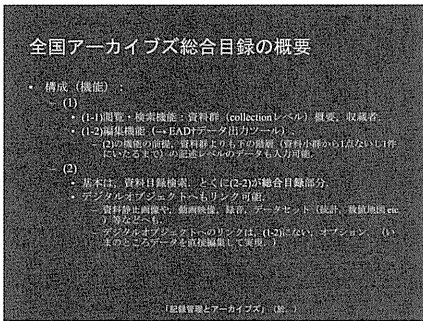
(B) 記述詳細



(B) 3) = 全国アーカイブズ総合目録

- 資料目録からデジタル画像まで、
- 運営：(現在) 国文学研究資料館。
- 構成：
 - 全国アーカイブズ総合目録
 - (1) 「東洋情報共有化データベース」
 - (1-1) 資料詳細閲覧・検索
 - (1-2) 資料情報編集・EAD/XML出力
 - (2) 検索システム
 - (2-1) 資料収集者ごと検索・閲覧
 - (2-2) 全資料収集者間検索・閲覧

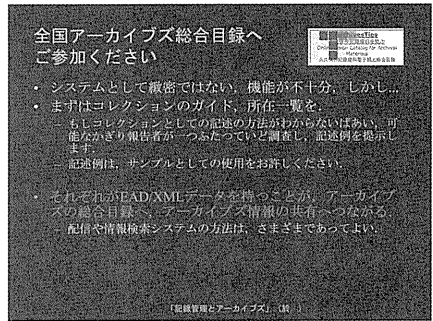
「記録管理とアーカイブズ」(18)



全国アーカイブズ総合目録の概要

- 構成(機能)：
 - (1) 閲覧・検索機能：資料館(collectionレベル)の閲覧・取業者。
 - (1-1) 編集機能：(EAD)データ出力ツール。
 - (1-2) 編集の前後、資料館レベルでの編集(資料館からできない1件1件にわたるまでの)記録レベルのアクセスも入力可能。
 - (2) 基本は、資料目録検索とくじりにの総合目録部分。
 - デジタルオブジェクトも利用可能。
 - 資料館と連携、編集機能、検索機能(統計、統計、編集機能 etc.) 等々あり。
 - アクセスプロトコルのリンクは、(EAD)ない、オブジェクト(1)を中心としたアクセス可能(検索)。

「記録管理とアーカイブズ」(18)



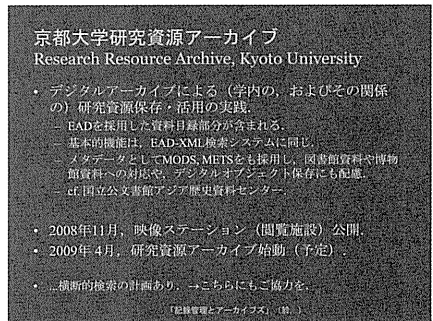
全国アーカイブズ総合目録へご参加ください

- システムとして厳密ではない、機能が不十分、しかし...
- まずはコレクションのガイド、所在一覧を。もしコレクションとしての記述の方法がわからない場合は、可能な限り報告がつまふつと調査し、記述例を提示します。
- 記述例は、サンプルとしての使用をお願いします。
- それぞれがEAD/XMLデータを持つことが、アーカイブズの総合目録とアーカイブズ情報の共有への鍵になる。
- 配信や情報検索システムの方は、さまざまであってよい。

「記録管理とアーカイブズ」(18)



京都大学 研究資源アーカイブ



京都大学研究資源アーカイブ Research Resource Archive, Kyoto University

- デジタルアーカイブによる(学内の、およびその関係の)研究資源保存・活用の実現。
 - EADを採用した資料目録部分が実装される。
 - 基本的機能は、EAD/XML検索システムに同じ。
 - メタデータとしてMODS、METSとも採用し、図書館資料や博物館資料への対応で、デジタルオブジェクト保存にも配慮。
 - cf. 国立公文書館アジア歴史資料センター
- 2008年11月、映像ステーション(開始施設)公開。
- 2009年4月、研究資源アーカイブ公開(予定)。
- 継続的検索の計画あり、こちらにもご協力。

「記録管理とアーカイブズ」(18)

参考文献・謝辞

- ・ 五島敏芳, 「日本におけるアーカイブズのオンライン総合目録構築において」, 『記録と史料』 18, 2008年3月, pp.1-17.
- ・ 本報告は, 国文学研究資料館史料館におけるアーカイブズの情報検索システムに関する研究と, EADを日本に導入する個人研究がもとになっています。ここまでたどりつくのには, 実はずいぶん多くの人たちにお世話になりました。ありがとうございます。
- ・ 報告の機会を下さり下さった国文学研究資料館・大学共同利用機関アーカイブズプロジェクトおよびKEK史料館に感謝いたします。

『記録管理とアーカイブズ』(録)