

デジタル・シルクロード:

多彩な文化遺産を統合するデジタルアーカイブ

北本朝展 大西磨希子 池崎友博 村松賢子 ドミニク・デュフ
マイヤー恵加 佐藤園子 エルハム・アンドロディ 山本毅雄 小野欽司
国立情報学研究所

「デジタル・シルクロード」プロジェクトは、シルクロード地域を対象とした文化遺産デジタルアーカイブの構築を目指すプロジェクトである。このプロジェクトの特徴は、情報科学と人文科学の学問的背景をもつメンバーが、お互いの強みを活かしながら共同作業を進めることにより、先端的なシステムと充実したコンテンツの両者をバランスよく実現しているところにある。またデジタルアーカイブを過去の文化遺産の保存という狭い枠組みにとどめず、コレクションからエグジビションへの組み換えや、現在進行中の事象のアーカイブ、さらには双方向性インタラクションに基づき成長するアーカイブなど、新しいアプローチにも実験的に取り組んでいる。本論文はこれらの試みを概括するとともに、その背景となる考え方と技術について述べる。

Digital Silk Roads:

A Digital Archive for Linking the Variety of Cultural Heritages

Asanobu Kitamoto, Makiko Onishi, Tomohiro Ikezaki, Takako Muramatsu, Dominique Deuff,
Eka Myer, Sonoko Sato, Elham Andaroodi, Takeo Yamamoto, and Kinji Ono
National Institute of Informatics

“Digital Silk Roads” is a project for establishing the digital archive of cultural heritage along the Silk Road. The project consists of members whose academic background is either computer science or humanities, so collaboration between members with different background contribute to the balanced evolution of the digital archive between advanced systems and rich contents. We also try to broaden the framework of traditional digital archives, which just aims at preserving existing cultural heritage, into new approaches, such as the recombination of collections into exhibitions, the digital archive of ongoing events, and evolvable archives based on bidirectional interaction. This paper summarizes these trials, and describes ideas and technologies behind new approaches.

1. はじめに

デジタル・シルクロード[1]は、シルクロード地域を対象とした文化遺産デジタルアーカイブの構築を目指すプロジェクトである。国立情報学研究所を中心に、ユネスコや他の学術機関・国際機関との共同研究を広げていくことにより、シルクロードの多彩な文化遺産のデジタルアーカイブを構築する計画を進めている。本プロジェクトの特徴は、情報科学と人文科学の学問的背景をもつメンバーが、お互いの強みを活かしながら共同作業を進めることにより、先端的なシステムと充実したコンテンツの両者をバランスよく実現しているところにある。またデジタルアーカイブの対象も多彩であり、シルクロードに関連する書

籍の電子化という研究史料デジタルアーカイブをはじめとして、自然災害やテロリズムによって危機に瀕する文化遺産を対象とした現在進行中の事象のデジタルアーカイブなど、デジタルアーカイブの枠を広げるための実験的な試みもおこなっている。さらに、デジタルアーカイブは単なる情報の保存庫ではなく多様な意味を生み出す場であるとの認識に基づき、リソースを多様に組み合わせる新たな文脈を生成するエグジビションといったアイデアも試みている。このように「デジタル」であることを活かした試みをおこなっていくことも重要な研究課題である。

本論文は、第2章でデジタル・シルクロード・プロジェクトの構成を紹介し、第3章ではデジタルアーカイブの基本設計について述べる。第4章で紹介す

るのは、現在ウェブでアクセスすることができる主要なコレクションである。続いて第5章ではコレクションを結合するための統合情報空間について説明し、第6章ではそれを拡張したエグジビションという概念を提案する。さらに第7章ではソフトウェアツールについて簡単に紹介する。第8章は現状の問題点と将来に向けての研究課題の総括であり、最後に第9章で本論文をまとめる。なお本プロジェクトはウェブサイト研究プラットフォームとしており、多くの研究成果はウェブサイト <http://dsr.nii.ac.jp/> で閲覧できる。

2. デジタル・シルクロード・プロジェクトの構成

本論文ではデジタル・シルクロード・プロジェクトの中で、以下のコレクションおよびツールについて論じるが、紙面の都合により、デジタル・シルクロード・プロジェクトの中でも省略したものがある。その詳細については文献[1]を参照していただきたい。それぞれの構成要素の2005年11月現在の進行状況については、以下に「公開中」や「未公開」といった形で記している。

2.1 コレクション

- 貴重書(「東洋文庫所蔵」図像史料マルチメディアデータベース)(公開中)
- 危機に瀕する文化遺産(イラン・バムの城塞、パーミヤーン大仏)(公開中)
- 写真(写真でつなぐシルクロード)(公開中)
- 日本人研究者論文集(未公開)

2.2 利用者別インタフェース

- 貴重書で綴るシルクロード(シルクロードに強く興味をもつ一般向け)(公開中)
- デジタル・シルクロード・キッズ(小学生から高校生、一般向け)(近日公開)

2.3 アーカイブ横断のための統合情報空間

- 地図(地理空間における統合)(未公開)
- 年表(時間空間における統合)(未公開)
- 用語(概念空間における統合)(未公開)

2.4 ソフトウェアツール

- MASS(多言語辞書管理ツール)(完成)
- RAMM(RDF アノテーション管理ツール)(開発中)
- DSR-CMS(コンテンツ管理システム)(完成)
- デジタル漂白処理(貴重書可読画像)(完成)

これらの要素の関係を表したのが図1である。まずウェブサイト全体はDSR-CMSというXMLに基づくコンテンツ管理システム(CMS)によって構築されている。また複数のコレクションが存在しており、それらを統

合情報空間が結合している。そしていくつかのソフトウェアツールが特定のコレクションのために開発されている。この未完成の部分も含めて、本論文ではプロジェクトの全体像を述べていくことにする。

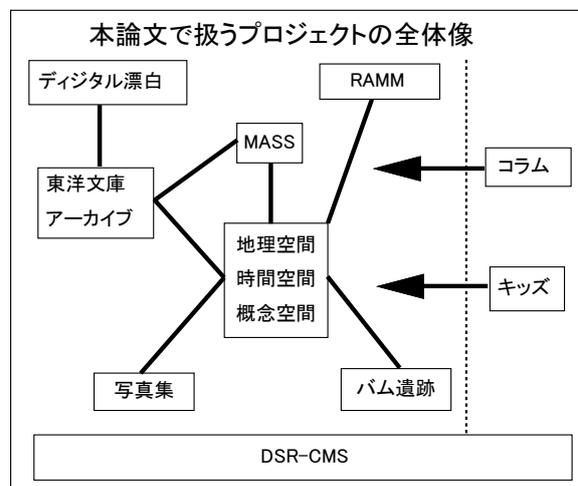


図1 本論文で扱うプロジェクトの全体像。

3. デジタルアーカイブの基本設計

3.1 コレクションの集合体

多彩な文化遺産のデジタルアーカイブを単一のシステムとして設計すれば、文化遺産の多様性を犠牲にする可能性がある。そこで本プロジェクトでは、まず文化遺産の対象ごとに自立したコレクションを構成し、それらを緩やかに結合した集合体として全体のデジタルアーカイブを完成させるアプローチを採る。現段階で公開が始まっているコレクションは、貴重書のコレクション、危機に瀕する文化遺産のコレクション、そして写真のコレクションである。これらは独立したサイト構成とナビゲーション、検索機能を備えており、それ自体がミニデジタルアーカイブとなっている。そして、これらを統合するための情報空間として地理空間や時間空間、概念空間などを導入する計画であるが、これについては第5章で述べる。

3.2 ナビゲーション指向のインタフェース

本プロジェクトでは、検索とナビゲーションという2つのインタフェースを併用することで、発見したいリソースへ到達するための手段を提供する。従来のデジタルアーカイブでは、複雑なキーワード検索機能を中心としたキーワード検索指向のインタフェースを持つものが多いが、これにはいくつかの問題点がある。第一に、検索すべきキーワードを明確に持たない一般の利用者は、適切なキーワードを入力しないと、中を覗き見することさえできない。第二に、利用者システムとキーワード体系が食い違っていると、たとえ関連キーワードを知っていても適切に検索することができない。第三に、「著者」などの適切なメタデータ欄にキーワードを入力しないと、思うような検索

結果を得ることができないが、これは現代の標準的な検索エンジンのインタフェースと比べて必要以上に複雑すぎる。こうした問題点を回避するために、本プロジェクトで構築するデジタルアーカイブは、ナビゲーション指向のインタフェースとなっている。すなわち、利用者の選択肢をできるだけ多数の明示的なリンクとして表示することにより、キーワードを入力するなどの利用者負担を軽減するものである。もちろん本サイトでも簡単な全文検索機能は提供しており、利用者が十分にシステムに慣れ親しんだ段階ではそれらを使いこなすことも可能である。しかし検索だけでは特に初心者にとって不十分であるというのが、本サイトの設計にあたって留意した点である。

3.3 多言語環境

シルクロードは世界中の研究者が研究対象とする地域であるため、研究史料も多言語にまたがる。また利用者は、日本国内のみならずシルクロード周辺諸国などにも散在している。そこで本サイトは多言語環境を基本とし、コンテンツとしてはヨーロッパ言語を中心に中国語なども扱い、それらを日本語・英語の2ヶ国語で提供することとする。またコレクションによっては、日本語・英語・ペルシャ語の3ヶ国語で提供するものもある。さらに言語横断検索や多言語で専門用語を翻訳するために、国立情報学研究所で多言語用語管理システム MASS (Multilingual-term Accumulation Support System) [1]も開発している。

3.4 メタデータ

シルクロードの文化遺産の多彩さに対応するためには、メタデータの管理にも柔軟さが必要である。多種のメディアに対応するだけでなく、それぞれの専門分野における語彙やその階層構造・関連性なども管理しなければならない。そこで国立情報学研究所ではメタデータ管理ツール RAMM (RDF Annotation Manager for Multimedia)を開発している。このツールについては第7章で述べる。

4. 主なコレクション

以下ではデジタル・シルクロードのコレクションとしてリソースの集積が進んでいるものを紹介する。

4.1 「東洋文庫所蔵」図像史料マルチメディアデータベース

<http://dsr.nii.ac.jp/toyobunko/>

シルクロードに関する基本的な文献を対象としたデジタルアーカイブであり、財団法人東洋文庫が所蔵する貴重書 40 冊、約 11000 ページの画像を提供している(2006 年 1 月までに 55 冊、約 14000 ページに増加する予定)。対象とする文献は著作権が消滅したものに限定している。表紙から裏表紙まで書籍の全ページを、高精細デジタルカメラ(一辺が 4000~5000 画素)を用いて歪みが少ない状態で撮影することで、地図や絵画などの詳細を閲覧するの

に十分な品質を確保している。ただしウェブサイトで公開している画像は低解像度に変換したものである。このアーカイブを構築するために用いたのは以下の技術である。

1. OCR(市販品)を用いた電子テキストの生成
 2. 形態素解析 ChaSen¹と検索エンジンGETA²を活用した電子テキストの全文検索
 3. 多言語管理ツール MASS を用いた類義語関係をたどる多言語検索
 4. デジタル漂白処理によるテキスト可読画像
- 上記の項目については文献[1]が詳しいが、そこで述べていない項目 4 については第7章で詳述する。

4.2 バムの城塞:地震を越えて残す記録

<http://dsr.nii.ac.jp/bam/>

2003 年 12 月 26 日に発生したイラン南東部地震で全面的に崩壊したイラン・バム遺跡の惨状を目の当たりにして、地震発生後数日後に開始した緊急デジタルアーカイブプロジェクトである。ただし開始時点ですでに遺跡は崩壊しているため、オリジナルの姿を正確にアーカイブすることは不可能な状況であった。そのため本プロジェクトでは、地震発生前の記録をできるだけ幅広く収集するために、ウェブサイトを通じて全世界に写真や映像の提供を呼びかけた。その結果、世界各地から 100 件以上の写真・映像資料が提供され、撮影時期・場所もバラエティに富んだ写真・映像コレクションを構築することができた。また高解像度衛星 QuickBird が地震発生前に撮影した画像を米国 Digital Globe 社(配給:日立ソフト)から入手することで、地震発生前の地理情報を正確に把握することも可能となった。現在の研究課題は、このように収集した不完全なデータを組み合わせ、地震前の姿をできるだけ忠実に再現するという課題にある。デジタルアーカイブは過去の文化遺産を記録するだけでなく、緊急災害対応のように現在進行中の事象もアーカイブすべきであると考え、これはその一例ともなっている。

4.3 写真でつなぐシルクロード

<http://dsr.nii.ac.jp/photograph/>

現代のシルクロードを撮影した写真のコレクションであり、過去に出版された貴重書に含まれるテキストやスケッチ、写真などと、現代の状況を比較することを目的とするデジタルアーカイブである。例えば100年前に撮影した写真と現代の写真と比較すれば、風景や風俗の変化などは容易に観察することができる。また上記の「バム遺跡」のように文化遺産そのものが消滅してしまえば、現代の写真は確実に後世への貴重な記録となる。つまり、現代に撮影した写真をデジタルアーカイブしていくことは、現在だけでなく将来にわたって、より大きな価値を生み出していく可能性が高い。このような観点から写真資料の収集もさらに拡大していく計画である。

¹ <http://chasen.naist.jp/>

² <http://geta.ex.nii.ac.jp/>

5. ハブとしての統合情報空間

このように、それぞれ独立したコレクションとしてデジタルアーカイブを構築した後に、これらを緩やかに結合するための統合情報空間を設定する。本プロジェクトで設定する統合情報空間は、以下の3つである。

- 地図(地理空間における統合)
- 年表(時間空間における統合)
- 用語(概念空間における統合)

この中で地理空間や時間空間は、実世界に対応する空間があるので想像しやすい。一方、概念空間は仮想的なものではあるが、用語として表現可能な概念の間の関係があらかじめ定義されている空間であると考えればよい。

このような情報空間を用いてリソースを結合するためには、まず個々のリソース(情報資源)を空間中の点または領域に射影する必要がある。例えば地理空間の場合は、地名というリソースを緯度・経度という形に変換して地理空間に射影する。また地域であれば、広がりをもつ領域を緯度経度領域に変換して射影する。したがって、あるテキストを地理空間に射影する場合には、出現する地名を順に地理空間に射影していけばよい。この作業は地名辞書があればより簡便である。次に時間空間であれば、イベントというリソースを期間という形に変換して時間空間に射影する。この作業は時代と年代とを対照した辞書などがあればより簡便である。最後に概念空間では、用語というリソースを概念空間の用語に射影する。この場合は完全一致する用語に射影すればよいが、場合によっては複数の関連用語に射影することが有効な場合もある。こうしてテキストに出現する用語を、順に概念空間に射影していけばよい。

このように個々のリソースを射影すると、今度は同じ点や領域に射影されたリソースの集合、あるいは近傍点や近傍領域に射影されたリソースの集合を取り出すことができる。つまり、こうしたリソースの集合を想定すれば、その要素となっている複数のリソースを結合することが可能となる。このように、統合情報空間をハブ(Hub)として複数のコレクションを結合する方法により、デジタルアーカイブを見通しよく構築していくことを考えている。

ソフトウェアとしては、まず地理空間での統合には地理情報システム(GIS)を導入する。本プロジェクトではすでにMapserver¹を導入して試験的に作動させている。また概念空間での統合にはMediaWiki²あるいは他のWikiクローンを導入する。この種のシステムは、用語がそのままURI (Uniform Resource Identifier)となるために空間を共有しやすいことが利点である。またWikiは記述を加えることで用語集としても活用できる。現状のMediaWikiは類義語関係などを陽に扱うことはできないが、Semantic MediaWikiなど、RDFを用いた用語関係の定義を試みる研究[2]もあり、今後の発展が見込まれる。

¹ <http://mapserver.gis.umn.edu/>

² <http://www.mediawiki.org/>

6. コレクションからエグジビションへ

6.1 概要

統合情報空間を用いたリソースの結合は、主に「点」としての結合、例えばある地名による結合やある時代、ある主題による結合という、局所的な結合に限られてしまう。このようなボトムアップな結合は限定された利用においては有用であるが、より広いテーマを見たい場合などには、全体像が見えないという意味で必ずしも便利ではない。もちろん統合情報空間全体を一望すればリソースの分布は可視化できるが、連鎖的な結合などは見えにくいという欠点が残る。となれば、俯瞰的な視点から複数リソース間の関係を明示的に可視化する方法を考える必要がある。そこで本プロジェクトでは、コレクションの再編成、すなわちデジタルアーカイブから「発掘」したリソースを新しい文脈のもとで並べ替えることにより、リソースの関係をより大局的に具体化する方法を考える。その模式図を図2に示す。

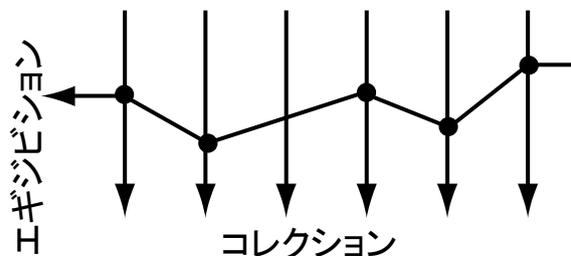


図2 コレクションからエグジビションへ。

個々のリソースはコレクションの文脈にしたがって並べられているが、コレクションを横断的に眺めてみれば、関連するリソースを横断的に見つけることができる。これらを結合するために、リソースをオリジナルの文脈から切り離し、新しい文脈の元に置いてみる。するとリソース間に新たな文脈を形成することができる。これはエグジビション(展示)の考え方に近い。

やや比喩的な言い方をすれば、あたかも既存の遺伝子の組み合わせを換えることによって新たな個性を持つ個体が誕生するように、リソースをもとの文脈から切り離して別のリソースと組み換えることにより、新たな文脈を生み出すことができる。そこで本論文ではこれを組み換え(Recombination)[3]とよび、このような組み換えによってデジタルアーカイブをより多面的に利用するための方法を追究する。

例えば、デジタルアーカイブのアクセス性を改善するために、利用者にて化したインタフェースを構築することも組み換えの一例として考えてみる。一般向きのインタフェース、子供向きのインタフェースなどを構築する問題は、デジタルアーカイブの中から利用者にとって適切なリソースを掘り起こし、それを適切に並べかえる問題と捉えることが可能である。

ただし、こうして生成された新しい文脈では誤解が生じる可能性があるため、常にオリジナルの文脈へ戻るためのリンクを設定しておくべきである。そしてそれは同時に、デジタルアーカイブへの新たな入口ともなる。こうして、いろいろな文脈に入口をもつ、入り込んでいきやすいデジタルアーカイブを実現することができる。また、作られた文脈は別の文脈の基盤として積み重ねていけるため、文脈を作れば作るほど、基盤となる材料が増えていくことにもなる。

6.2 貴重書で綴るシルクロード

<http://dsr.nii.ac.jp/rarebook/>

「東洋文庫アーカイブ」を中心として写真コレクションなどからもリソースを発掘し、それをストーリーに沿って並べ替えて提示する構成のエグジビションである。このエグジビションは、以下の手順で構築している。

1. 専門的知識をもつキュレータが、一つのテーマを設定する。
2. テーマに関連すると考えられるリソースを、リソースの一覧表示や全文検索などを組み合わせ、デジタルアーカイブの中から発掘する。
3. テーマに沿った一つのストーリーを考案し、プロットに合わせて発掘したリソースを並べ替える。
4. 不足したリソースがあればデジタルアーカイブの外で調査を加える。
5. リソースとストーリーを整えてウェブで公開する。こうして配置されたリソースは常にオリジナルの文献とリンクされているため、オリジナルの文献での文脈とキュレータが設定した文脈の両方を参照しながら、シルクロードに関するストーリーを読み進めていくことができる。こうして、より大局的な観点からデジタルアーカイブを閲覧していくことが可能となる。

ただし、こうした試みをさらに大規模化するためには、キュレータのコミュニティが必要である。なぜならば、デジタルアーカイブの多彩な文化遺産を十分に発掘していくためには、それだけ多彩な専門家の目が必要となるからである。今後は、多人数のキュレータによるデジタルアーカイブ発掘作業が、エグジビションの発展には必要になると考えている。

6.3 デジタル・シルクロード・キッズ

「貴重書で綴るシルクロード」はシルクロードに強い関心をもつ一般の人々をターゲットとしたインタフェースである。しかしより親しみやすい子供向けや一般向けのインタフェースも多くの人にとっては有用なものである。これも、デジタルアーカイブから子供に適したリソースを発掘して並べるという意味ではエグジビションの一種であるが、この場合はリソースをそのまま提供するのではなく親しみやすい形に改変して提供することが重要な課題となる。例えば貴重書のスケッチや写真などを参考に、ぬり絵を描きながらシルクロードに対する理解を深めていくプログラムや、シルクロード百科として子供向けに用意した用語集からシルクロードを学ぶプログラムなど、教育や学習に有用なプログラムを準備していく計画である。

7. ソフトウェアツール

本章ではデジタル・シルクロードで用いているソフトウェアツールのうち、本プロジェクトで開発したものの中から特に、(1) RDF アノテーション管理ツール、(2) デジタル漂白処理、の二点を取り上げて述べる。

7.1 RDF アノテーション管理ツール

デジタルアーカイブにおいては、リソースに対するメタデータを管理することが重要な仕事のひとつである。このようなメタデータについてはISAD(G)やEADなどアーカイブを対象としたメタデータや、Dublin Coreなど書誌を対象としたメタデータなどが国際標準として定められている[4]。しかし実際にメディアの中身にまで立ち入ってメタデータを記述しようとすれば、これら国際標準のみでは限界がある。そこで本プロジェクトで構築するRAMM (RDF Annotation Manager for Multimedia) では、RDF (Resource Description Framework)¹を用いてメタデータを記述することにより、RDFの記述力を活用したメタデータの管理と検索を実現する。リソースのメタデータは、現在のところ以下の3層に分解して記述するモデルを用いている。

概念層	概念に対応するメタデータ
実体層	実在物に対応するメタデータ
媒体層	ファイルなどに対応するメタデータ

図3 リソースのメタデータの3階層モデル。

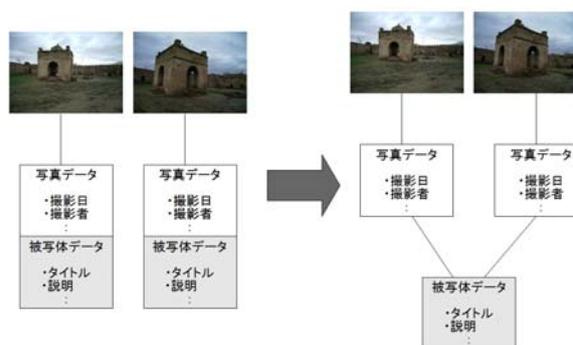


図4 媒体層と実体層との分離の例。

例えば図4のように、画像ファイル(写真)にメタデータを付与する場合を考える。ここで必要なのは、撮影日や撮影者など写真自体のメタデータと、被写体そのものについてのメタデータである。被写体が同一の複数の写真の間では、被写体そのもののメタデータは共有すべきであるから、逆に写真自体のメタデ

¹ <http://www.w3.org/RDF/>

ータとは分離すべきであると考えられる。したがって本システムでは、ファイルなどについてのメタデータを媒体層で記述するとともに、実世界のあらゆるもの(実在物)を個々にリソースとして定義し、これらのメタデータを実体層で記述する。そして実体層の概念に関するメタデータを概念層で記述する。

概念層では概念間の抽象・具象関係に限定した階層構造の体系を用いる。また複数の分類基準を許容するために、一つの実在物リソースが複数の概念に属せるようにする。また、概念ごとに、そこに属するリソースの属性セットを定義可能とし、事前に定義された基本概念と基本属性に加え、利用分野に応じた専門的な属性を利用者が定義するための機能も付け加えている。

ただし RDF は自由度が高いため、記述に任意性が生じて記述が不統一となる可能性がある。例えば図5のように実在物がある実在物の一例である場合、実体層のリソースを連結した構造が一つの入力方法として想定されるが、これに一意に定まるわけではないので、利用者は別の方法で入力する可能性がある。こうした問題は、入力ガイドラインの作成とユーザインタフェースの工夫によって部分的には回避できようが、根本的にはメタデータの構造にあいまいさが少なくなるように検討を加えていく必要がある。

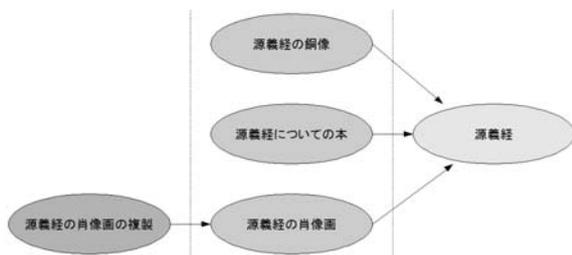


図5 実体層の分解



図6 部分画像へのメタデータ付与

さて、本システムの特徴的な機能の一つに、部分画像へのメタデータ付与という機能がある。例えば図6のように、画像という一つのファイルの中に、いくつかの実在物が含まれているとしよう。このような場合には、全体画像にメタデータを付与するよりも、画像の

部分に対してそれぞれメタデータを付与できれば、より正確で意図が明確なメタデータの記述が可能になる。本システムでは、四角形で囲むなどの方法で部分画像を定義し、それぞれにメタデータを付与することができるようになっている。

こうして入力したメタデータの検索については、単純な全文検索に加えて、属性を用いた検索もサポートしている。また概念を用いた検索では、上位概念を指定することで下位概念のリソースも検索することができる。概念が検索できれば、そこに属するリソースとリンクする画像や音声などの媒体層ファイルを一覧する機能もある。この機能と、部分画像へのメタデータ付与機能とを組み合わせると、例えばある特別な装飾模様を検索したときに、壁画の一部に含まれる模様と、陶器の絵付に含まれる模様を一括して出力することも可能である。そして部分画像から親画像へとたどることで、全体の中での位置も確認できる。

今後の課題としては、まずフレームワークとして用いている RDF の能力をフルに活用できるより複雑な検索、例えば「この絵を描いた人物の一族全員の作品は？」といった検索を可能とするように、ユーザインタフェースを拡張するという課題がある。また概念層の実装をよりシンプルで整合性の取れたものにも第二の課題である。第三の課題は、本システムの本来的目的である研究者の共同作業環境の実現のために、グループスペースの設置やメタデータの付与、承認、公開という一連のフローを管理するための機能、作成されたデータの公開機能などについても実装を進めていく計画である。

7.2 デジタル漂白処理

これは貴重書デジタルアーカイブのために開発した画像処理手法である。その目的は、本プロジェクトのような公開型デジタルアーカイブにおいて、複製可能性を不必要に高めることなく、貴重書のアーカイブで最も重要な要件であるテキスト可読性を高めることにある。本手法は以下の相反する二つの要求をうまくバランスさせることを目標とした手法である。

- 小さな文字のテキストも読めるように、できるだけ高い解像度で公開したい。
- 画像の不正利用を防ぐために、できるだけ低い解像度で公開したい。

後者の不正利用は、公開型デジタルアーカイブにおける最大の懸念材料である。従来の対策には電子透かしなどの方法があるが、あまり有効な対策とはなっていない。そこで本プロジェクトではデジタル漂白処理(digital bleaching)[5]を提案し、画像解像度と画像品質とをうまくバランスさせるという単純な考え方で、この問題を解決することを考える。

その基本的アイデアは以下のとおりである。貴重書のデジタルアーカイブにおいて、複製価値が高いのは図像、すなわちスケッチや絵画・写真であること

に着目すると、テキスト可読性の向上と同時に図像の複製価値が向上しなければ、当面の対策としては十分であることがわかる。そこでまず、テキスト可読性を向上させるために、画像の色情報や、劣化した紙のシミや汚れ、撮影方法に起因する微妙な陰影などを除去し、紙面領域を一律な背景として目立たなくさせることを考える。そのためには、紙面領域の代表値に対応するハイライト値と、文字領域の代表値に対応するシャドウ値とを画像から自動的に推定し、自動的にコントラスト調整をおこなえばよい。具体的には、紙面領域が空間的に一律な濃度値になるようにハイライト値を選択すれば、紙面は漂白した紙のように一律な色となって背景と認識されるようになり、この効果によって利用者は文字のみに集中してテキストを読むことができるようになる。こうして、テキスト可読性は向上するはずである。

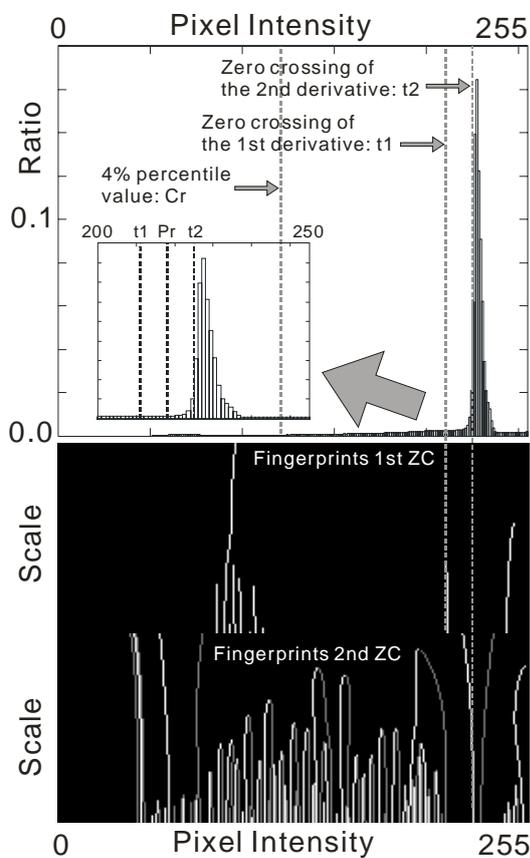


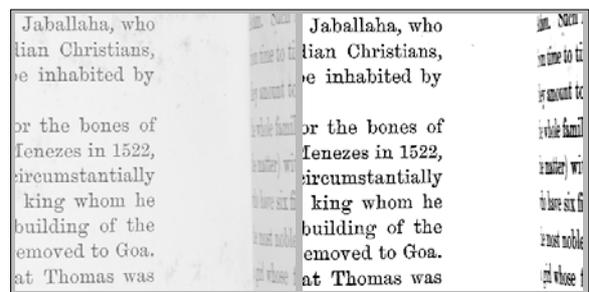
図7 スケールスペースヒストグラム法。

このような処理を実現するために、本論文ではスケールスペースヒストグラム法を用いる。この手法は、画像ヒストグラムという1次元信号の各種ゼロ交差 (Zero Crossing: ZC) を多重解像度の考え方で解析し、そのスケールスペース上での振る舞い (Fingerprint) を追跡することにより、構造的に顕著なゼロ交差を検出する方法である。ゼロ交差の位置はヒストグラムの形状的な特徴に対応しており、局所最小 (1次微分) ゼロ交差はヒストグラム形状の谷、2次微分ゼロ交差はヒストグラム形状の変曲点に対応す

る。この性質を用いてヒストグラムの構造的に顕著な特徴を検出することができる。紙面領域の代表値を推定する方法の概略は以下ようになる。

1. 紙面領域の画像中での割合や代表値の濃度などに関する確率密度関数を想定しておく。これは複数の候補がある場合のランキングに利用する。
2. 画像ヒストグラムの顕著な局所最小 (ゼロ交差) の中から、濃度が最大に近いものを t_1 とする。
3. 画像ヒストグラムの顕著な2次微分ゼロ交差の中から、 t_1 以上、かつ正から負に変化するゼロ交差の濃度を t_2 とする。
4. 紙面領域の代表値を $(t_1 + t_2) / 2$ とする。またゼロ交差の候補が複数ある場合には、この代表値を候補ごとに計算して 1 で用意した確率密度関数に代入し、出現確率が最大のもの (尤度が最大のもの) を最終的な代表値とする。

この手順を図示したものが図7である。また、こうして生成した漂白処理画像を、単純に白黒濃淡化した脱色処理画像と比較した結果を図8に示す。脱色処理画像では書籍を撮影した写真という印象を拭い去ることが難しいが、漂白処理によって紙面が一律な色になれば、紙に印刷した書籍に近い印象を与えることができる。またコントラストの強調により、テキストそのものも脱色処理画像より読みやすくなっている。なお「東洋文庫所蔵」図像史料マルチメディアデータベースでは、この漂白処理を約 11000 ページの画像に対して全自動でおこなっている。そのため提案手法であるデジタル漂白処理は、ある程度の書籍の変動には頑健であるといえ、他の貴重書画像に対してもロバストに適用できると想像している。



脱色処理画像 漂白処理画像

図8 脱色処理画像と漂白処理画像の比較。

8. 議論

8.1 現状の問題点

以上述べてきたように、デジタル・シルクロード・プロジェクトでは多彩な対象を選び、情報学的手法を適用しながらデジタルアーカイブへの新しいアプローチを続けてきた。それに対して、利用者側から指摘されている問題点をまとめると、以下ようになる。

1. コンピュータを扱うスキルが高くない利用者への対応が不十分である。このアーカイブの主な利用者層の一つとして想定している東洋史や

美術史の研究者の中には、あまりコンピュータの利用に慣れていない人もいるので、使い方のガイドが必要ではないか。

2. 専門用語が頻出するために、専門用語の解説がないと使いこなすことが難しい。原資料だけではなく、専門用語のレファレンスなどが必要ではないか。
3. このアーカイブを利用して発見した情報をアーカイブにフィードバックし、他の人と共有するための手段がない。せっかく発見した情報も今のままでは埋もれるだけであるので、例えば付箋のような形で情報を共有できるようにしてほしい。

8.2. 成長するデジタルアーカイブを目指して

特に最後の指摘は、デジタルアーカイブの双方向性にまつわる重要な問題提起であり、今後のデジタルアーカイブの発展の方向性を示唆するものでもある。これまで本プロジェクトでは、研究者に対して研究素材を提供することを名目にデジタルアーカイブを構築してきた。しかしこれだけでは、研究素材を集積しているデジタルアーカイブという「場」を有効活用しているとは言えない。せっかく研究者がデジタルアーカイブという場を共有しているのに、現状ではその場を共有する手段がない。もし研究者が何らかの情報を残せるようにし、それを他の人も共有できるようにすれば、デジタルアーカイブが情報共有の場へと発展していくことが可能なのではないか。

このように、研究者が得た情報をうまくフィードバックできるようなメカニズムが実現すれば、フィードバックされた情報が集積すればするほど利便性が向上していき、やがて自律的に成長するデジタルアーカイブを実現することも可能になるだろう。このような研究者支援のメカニズムをどのように実現するかが、今後の大きな研究課題である。

デジタルアーカイブは原資料を忠実に保存する場であるという狭義のデジタルアーカイブ像からすれば、こうした個人的な情報を原資料と共にアーカイブするものは、アーカイブの範囲を逸脱しているかもしれない。しかし私は、このような成長するデジタルアーカイブのモデルが中国の史書にあると考えている。

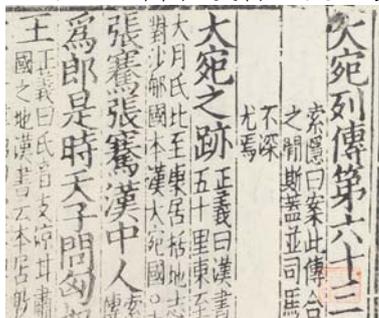


図9 史記 第二百二十三卷 大宛列傳 第六十三

図9に示すように中国の史書は、オリジナルのテキスト(1段組)を記録するだけでなく、後世の学者が加えた注釈(2段組)も一体のテキストとして記録して

いる。つまり、価値ある注釈は、原資料と一体となって書籍にアーカイブされている。この精神を受け継ぐならば、原資料の忠実な保存に加えて、研究者による注釈なども含めてアーカイブを構築していくことにも価値があると考えられる。もちろん中国の史書でもそうであるように、原資料と付加情報とはフォーマットの上で区別されるべきである。しかし、個人がつけた注釈が整理されれば、それも一つの知的財産として後世に受け継がれていくことが期待できるだろう。

9. おわりに

本論文では多彩な文化遺産を統合するデジタルアーカイブの例として、デジタル・シルクロード・プロジェクトの概要を説明した。本プロジェクトは、情報科学のメンバーがシステムの設計を担当し、人文科学のメンバーがコンテンツの充実を図ることで、両方の強みを活かしたデジタルアーカイブの構築を進めている。また研究プロジェクトとして実験的な試みにも積極的に取り組んでいる。第8章の議論でも述べたように、研究者コミュニティとの連携は今後の大きな課題である。また「バム遺跡」のプロジェクトである程度は成功したように、一般の人々と連携して広くデータを収集する方法も有効である。これらの課題をまとめると、双方向性を備えたデジタルアーカイブの構築が大きな課題であるとまとめることができる。

謝辞

貴重書の利用に関して、財団法人東洋文庫の斯波義信博士、田仲一成博士の多大なる協力に深謝する。なお本研究は、科学研究費補助金(研究成果公開促進費・データベース)の助成を受けている。

参考文献

- [1] Ono, K., Yamamoto, T., Kamiuchi, T., Andres, F., Kitamoto, A., Sato, S., Andaroodi, E., "Progress of the Digital Silk Roads Project," *Progress in Informatics*, Number 1, pp. 93-141, March 2005.
- [2] Muljadi, H., Takeda, H., Araki, J., Kawamoto, S., Mizuta, Y., Demiya, S., Suzuki, S., Kitamoto, A., Shirai, S., Ichiyoshi, N., Ito, T., Abe, T., Gojohori, T., Sugawara, H., Miyazaki, S., Fujiyama, A. "Semantic MediaWiki: A User-Oriented System for Integrated Content and Metadata Management System," *International Conference on WWW/Internet*, 2005.
- [3] Kitamoto, A., Sato, S., Yamamoto, T., Ono, K., "Context Recombination for Digital Cultural Archives," *Proceedings of the International Conference on Digital Archive Technologies (ICDAT2004)*, pp. 105-119, 2004.
- [4] 国文学研究資料館史料館編、アーカイブズの科学、柏書房、2003.
- [5] 北本 朝展, 山本 毅雄, 佐藤 園子, ナイジェル コリアー, 川添 愛, 小野 欽司, "貴重書デジタルアーカイブにおけるテキスト可読性と異種メディア間共参照アノテーション", *画像電子学会誌*, Vol. 33, No. 5, pp. 737-745, 2004.