

KU-ORCAS

——オープン・プラットフォームが切り拓く新しい人文知の未来

内 田 慶 市

はじめに

KU-ORCAS（関西大学アジア・オープン・リサーチセンター）は、2017年度文部科学省私立大学研究ブランディング事業に採択されたもので、その目的は関西大学の特色ある豊富なりソースを基盤とする東アジア文化研究のデジタルアーカイブを構築し、その活用を通じて東アジア文化研究の世界的ハブ的研究拠点としてのブランドを確立することにある。

関西大学の東アジア研究は今から約250年前の江戸時代の「泊園書院」に遡るが、その泊園書院を源とする東西学術研究所（1951年創設）を中心に展開され、特に、2005年以降、文科省の学術フロンティア推進事業による「関西大学アジア文化交流センター（CSAC）」（2005–2009年）、私立大学戦略的基盤形成事業による「関西大学アジア文化研究センター（CSACII）」（2009–2013年）、更には2007年から2011年までの文科省グローバルCOEの採択といった成果を挙げてきており、このKU-ORCASはそうした研究成果の蓄積の上に打ち立てられたものである。

1. CSAC デジタルアーカイブの現状

さて、まず最初にCSAC及びCSACIIの研究プロジェクトで構築した「CSAC デジタルアーカイブ」について簡単に述べておく。

1.1 近代漢語文献データベース（2006年より）

このデータベースは2006年に筆者の科研費によって構築した近代漢語文献

1.2 文献データベース

関西大学には下記のような東アジア関連の個人文庫が多数所蔵されている。このデジタル化にもこれまで鋭意取り組んできている。

- 内藤文庫 (33500点) ……内藤湖南 (漢籍)
- 長澤文庫 (30497点) ……長澤規矩也 (国漢籍)
- 中村文庫 (33491点) ……中村幸彦 (国文)
- 増田文庫 (16184点) ……増田渉 (魯迅、西学東漸)
- 吉田文庫 (2479点) ……吉田伊三郎 (アジア外交)
- 鬼洞文庫 (10309点) ……出口神暁 (国文)
- 泊園文庫 (16954点) ……藤澤東暎・南岳・黄鶴・黄坡



図3 関西大学デジタルアーカイブ (<https://www.iiif.ku-orcas.kansai-u.ac.jp>)

現在までに約6000冊程度のデジタル化が完了しているが、公開しているのはそのうち3000冊程度である。

なお、こうした個人文庫には書籍以外にも、書簡類、書画類や非文献資料も多数所蔵されており、そうした資料のデジタル化・公開も行っている。

例えば、以下のようなものがある。

漢封泥のデジタル化（20件程度）

内藤湖南の書簡が8000件程度：撮影終了

泊園関係の印鑑類が300点程度

内藤湖南文庫所蔵の軸物や貴重書庫収蔵の軸物700点程度（公開待ち）

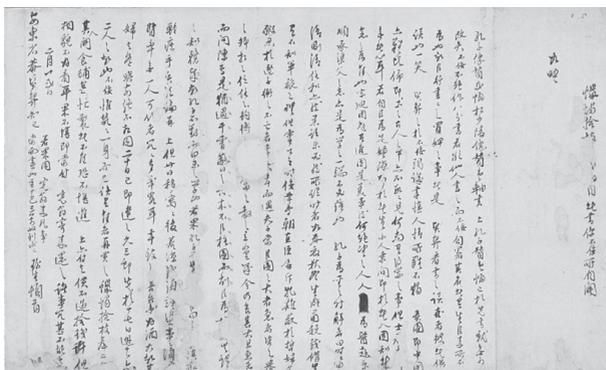


図4 朱舜水先生手簡（内藤文庫）



図5 羅叔言參事
臨秦權條幅
（内藤文庫）

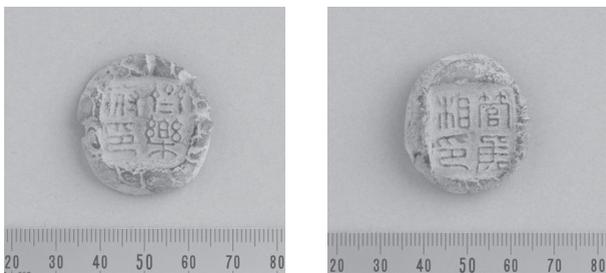


図6 漢代の封泥（Chinese-Style Wax Stamps）

2. アジアにおけるデジタル化の現状

ところで、アジアにおけるデジタル化の現状であるが、最も早く（約20年前）から、しかも大規模にデジタル化を行ってきたのは、やはり台湾中央研究院の歴史語言研究所漢籍電子文獻資料庫であろう。現在、約5億字のデジタル化が行われており、全文検索も可能である。

また、中国でも盛んに行われてきており、例えば、CADAL (China Academic Digital Associative Library=中国数字図書館国際合作計画) は浙江大学と中国工程院による国家的プロジェクトとして2001年に開始され、すでに700万冊のデジタル化が完了している。ただ、日本ではこれまで国会図書館や私どもの関西大学も連携して加入していこうと考えてはきたが実現には至らなかった。最近ようやく東京大学が日本で初めて参画している。

この他、環太平洋デジタル図書館連合 (PRRLA=Pacific Rim Research Library Alliance) という国際的な組織もあり、香港、中国、台湾、マカオ、オーストラリア、カナダ、アメリカなど33の大学が参加しているが、これも残念ながら日本の大学は加盟していない（関西大学は筆者が図書館長を務めていた時代に一度加盟したが、現在は幽霊会員となっている状況である）。

日本に目を向けると、早稲田大学や国会図書館近代ライブラリー等は早くか

Chinese Text Project

Welcome

Welcome to the Chinese Text Project homepage. The Chinese Text Project is an online open-access digital library that makes pre-modern Chinese texts available to readers and researchers all around the world. The site attempts to make use of the digital medium to explore new ways of interacting with these texts that are not possible in print. With over thirty thousand titles and more than five billion characters, the Chinese Text Project is also the largest database of pre-modern Chinese texts in existence.

You may wish to read more about the project, view the [pre-Qin and Han](#), [post-Han](#) or [Wiki](#) tables of contents, or consult the [instructions](#), [FAQ](#), or list of [tools](#). If you're looking for a particular Chinese text, you can [search for texts by title](#) across the main textual sections of the site.

網站有中、英文版本，也有繁、簡體版，可透過每頁左上角的連結隨時調整。

Date	Content
2016-10-10	Harvard Yenching Library Chinese materials added Thanks to the support of Harvard Yenching Library, over 5 million pages of scanned materials from the Yenching Library collection have been added to the Library section of the site, including high quality images from the Chinese Rare Books collection. Approximate transcriptions created using the ctext.org OCR procedure have also been added to the Wiki, making these materials full-text searchable. In future we hope to collaborate with other libraries to include materials from their Chinese language collections.
2015-07-02	Support for Unicode 8.0 adds new characters A new version of the Unicode standard has been released, defining thousands of additional rarely used and variant Chinese characters. Support for these has been added to the dictionary section of the site; to view these characters, please install the latest version of the Hanazono font . Many new characters belong to "CJK Extension E" - you can confirm system support for these from the Font Test Page .

Previous technical updates are listed in the "Latest updates" forum. Please note that only major technical updates are listed here; for content updates and recently added texts (which occur daily and hourly), please [log in](#) and refer to the [Wiki Recently Updated](#) sections.

図9 CTEXT (Chinese Text Project) のページ (<https://ctext.org>)

らデジタル化を推進してきたが、ここに来て、国文学研究資料館の日本語の歴史的典籍データベース構築プロジェクトや京都大学、島根大学等々も積極的にデジタル化を推し進めてきている。

一方世界における東アジア文献資料のデジタル化の状況について簡単に触れておくと、例えば以下のような機関で多くの文献のデジタルアーカイブが実現している。

Gallica (BnF=フランス国立図書館)

CrossAsia (ドイツ国家図書館)

Münchener DigitalisierungsZentrum Digitale Bibliothek (MDZ=ミュンヘンデジタルセンター=バイエルン州立図書館、IIIF 対応)

Digitalisierte Sammlungen (ベルリン州立図書館)

Digital Vatican Library (DVL=DigiVatLib)

イエズス会文書館 (Archivum Romanum Societatis Iesu)

Heidelberg University (COE)

Harvard-Yenching Library (ハーバード大学)

Serica (Bodleian Library, オックスフォード大学)

National Library of Australia (NLA=オーストラリア国立図書館)

もちろん、Hathi Trust、Internet Archive、Google books といったいわゆるオープンアクセスサイトも充実してきている。

この他、中国語に特化したものとして以下のようなものを挙げておく。

中央研究院近代史研究所「英華字典」

CTEXT (中国哲学書電子化計画)

CCL (北京大学中国語学中心)

BCC (大数据与語言教育研究所)

3. サイロ問題の打破

現在、国内外を問わず、多くの機関、組織でデジタル・アーカイブが進められてきており、デジタル化は今や世界の潮流であるが、ここに一つ大きな問題が存在する。それは「サイロ問題」と言われるものである。

つまり、多くの研究機関でそれぞれにデジタル・アーカイブ化が進められているが、横の連携が希薄で、各所で同じようなものが別々に作られて「閉じた」状態となっている。そして、自分たちのサイトに貴重な資料が沢山保管さ



図10 IIFの実例 (KU-ORCAS 所蔵の英華字典の対照比較)

れていても、閉じた状態であるため、アクセスも制限され、放っておく（アクセスしない）と腐ってしまう（アクセスが少ないからと維持されない）わけである。

こうしたサイロ問題が生まれる背景には、例えば、画像における国際的規格がないということがある。

こうした状況を打破するために、近年提唱されているのが、IIF（トリプル・アイ・エフ＝International Image Interoperability Framework 国際的画像相互運用の枠組み）という試みである。

このIIFとは国際的なWebコンテンツ共有の枠組みであり、国外ではスタンフォード大学図書館、英国図書館、フランス国立図書館、オックスフォード大学ボドリアン図書館等多くの図書館がすでにこれを採用しているが、日本でも東京大学大学院人文社会系研究科次世代人文学開発センター人文情報学拠点、京都大学図書館機構、関西大学アジア・オープン・リサーチセンター、国会図書館、国文学研究資料館等がこれを取り入れている。

このIIFのメリットとしては、同じビューワー内に他機関のIIF対応資料画像を同時に見ることが可能であり、アノテーション（いわゆるメモ、注釈）機能を追加することも可能となり、資料の比較研究が容易となるほか、多様な立場の人の様々な解釈を共有できるようになることが挙げられる。今後、日本でも多くの機関がこの規格に準拠したデジタル化を行うことが望まれる。

4. デジタル・アーカイブと人文研究

上ですでに述べたとおり、デジタル・アーカイブ化は昨今の人文学研究の一つの大きな潮流となっている。それは学問研究に便宜を提供するだけでなく、実は人類の「知の遺産」の「保存」ということにもつながっている。またこうしたデジタル化、公開という流れはまさに筆者が以前から主張してきた「秘蔵は私蔵なり」という主旨にも完全に合致するものである。あらゆる文献資料は「公開」されるべきである、それが「書の使命」であるというのが筆者の基本的考え方である。

ただ、単にデジタル化、アーカイブ化だけでは不十分である。それは、研究者の立場から言えば、新しい研究方法と結びつく必要があるのだと考えている。

例えば、かつての語彙研究と言えば、カードの枚数に比例するものであり、ある語彙の初出は読んだ量に依拠するものであった。それが今やコーパス、コンコーダンス、全文検索、全語彙索引といったものによって研究方法は一変した。

また、こうしたツールを利用することで単なる語彙史研究から他の研究領域（概念史研究、思想史研究等）との領域を越境したコラボも実現してきている。様々な情報を追加していくことで、様々な角度からのアプローチが可能になるのである。

例えば、Google の Ngram Viewer を利用すれば、1500年頃から現在までのトレンドキーワードを調べられる。

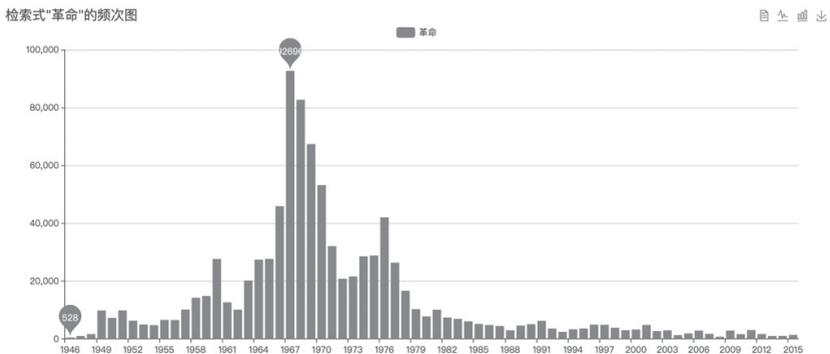


図11 Ngram Viewer での「革命」の検索結果

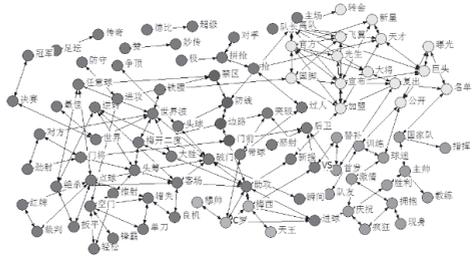
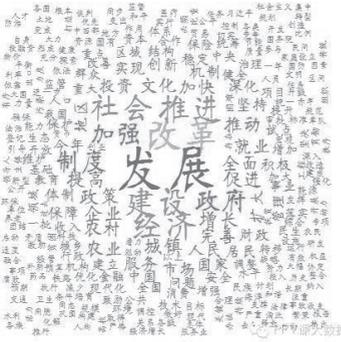


図12 中国語のテキストマイニングの例

Bibliographic Information	Version 表: 表華本	Version 九: 九江本
官話指南各種 Original Source 表華本 九江本 翻語指南本 Witness List <ul style="list-style-type: none"> Witness 表: 表華本 Witness 九: 九江本 Witness 翻: 翻語指南本 Electronic Edition Information: Publication Details: Publication Information 《官話指南》是明治時代日本人最早編的漢語課本。作者是吳松如大和館本所。《官話指南》初版於1867年。後來有各種版本。如上海語源、廣東語源等。請參照內田康市・水野善寬編《官話指南的語言研究》。	1 您納 貴 姓？ 魏姓吳。 請教台甫。 草字寶號。 寶 姓 姓 幾 位？ 我 們 2 前幾 天 我 去 的 時 候。 他 也 托 我 問 您 好 家 事。	1 您 貴 姓？ 魏姓吳。 請教台甫。 草字寶號。 寶 姓 姓 幾 位？ 我 們 2 前幾 天 我 去 的 時 候。 他 也 托 我 問 您 好 家 事。

図13 TEI を利用した『官話指南』の校勘例

テキストマイニングでは、文章のデータを単語や文節で区切り、それらの出現の頻度や共出現の相関、出現傾向、時系列などを解析することで有用な情報を取り出すことが可能である。

この他、TEI (Text Encoding Initiative, 1987-) という人文学研究のための電子テキストの効果的効率的な共有のためのガイドラインあるいはその構造化のルールセットを利用することでテキスト (XML 化) は汎用性、永続性を持つことになり、例えば、テキスト校勘などを視覚的に行うことができる。

5. 中国語とコンピュータの古くて新しい関係

中国（中国語）とコンピュータは実は極めて古くから関係がある。易の原理は陰陽の二元論であり、それはコンピュータの二進法に通じている。易の八卦は3ビット（三爻の組み合わせ）であり、上下で64通りの組み合わせ（すなわち6ビット）によって森羅万象を表現する。

ただし、古い関係にありながらも、いわゆる漢字や日本語はコンピュータに乗りにくい言語であった。私がコンピュータを始めた90年代でも「文華」等々のソフトがあったが、中国語の漢字を表示したり入力したりする場合は、第2水準の漢字ボードを中国語の漢字ボードに置き換えたりしていた。それが今ではユニコードの時代となり何万もの漢字を自由に扱うことができるようになっている。

それでも、中国語や日本語の自然処理は欧米語のように簡単にはいかないのだ。

今から約20年前に当時出入りしていたニフティサーバ（NIFTY-Serve）のFPRINT-15「せどおうくばある（SED, AWK, PERL）」のフォーラムで知り合った彌永信美と齋藤希史に依頼してPerlとApple Scriptで「全語彙索引作成プログラム」（1999.2.3）を作成したことがある。これは今でも利用できるし、現在は氷野善寛氏によりウェブ上でのプログラム（<http://www.chlang.org/contents/index-converter/>）が存在する。

使い方は、至って簡単で、まず以下のようなテキストを準備（文字コードはUTF-8）し、プログラムにかけるだけである。

[sample 1]

我 是 關西大學 的 學生 我 今=（これは次の行の単語と1語であるという印）
年 二十一 歲 我 學 漢語 專業 我
住在 大阪 我 有 父親 母親 弟弟
我 弟弟 今年 十七 歲 明年 要 考 大=
學 我 父親 今年 五十四 歲 母親 五=
十二歲

また、巻数、ページ数、行数を結果に表示させたいときには、先のテキストに次のように標識を付ける。

<V 1> ……巻数

<P 1a> ……ページ数（この例では、1葉の表。裏の場合はbで表す）

<L 1> ……本文の最初の行を1行と示す

結果は以下のようなになる。

idx result of the file

/Volumes/MyDocument/Users/ni/Desktop/test_things/idx.pl/idx_pl_Unicode_version/test_folder/chinese_test.txt

二十一：1 (1-1a-2)
五十二歳：1 (1-1a-6)
五十四：1 (1-1a-5)
今年：3 (1-1a-2, 1-1a-4, 1-1a-5)
住在：1 (1-1a-3)
十七：1 (1-1a-4)
大學：1 (1-1a-5)
大阪：1 (1-1a-3)
學：1 (1-1a-2)
學生：1 (1-1a-1)
專業：1 (1-1a-2)
弟弟：2 (1-1a-3, 1-1a-4)
我：7 (1-1a-1, 1-1a-1, 1-1a-2, 1-1a-2, 1-1a-3, 1-1a-4, 1-1a-5)
明年：1 (1-1a-4)
是：1 (1-1a-1)
有：1 (1-1a-3)
歳：3 (1-1a-2, 1-1a-4, 1-1a-5)
母親：2 (1-1a-3, 1-1a-5)
漢語：1 (1-1a-2)
父親：2 (1-1a-3, 1-1a-5)
的：1 (1-1a-1)
考：1 (1-1a-4)
要：1 (1-1a-4)
關西大學：1 (1-1a-1)
Total words：24

ただ、このテキスト入力で最大の難点が「単語の区切り」である。中国語や日本語はこの単語をどこで切るかが極めて問題となるのだ。間にスペースを入れていくのだが、これまでは手入力で行ってきた。

最近ようやく、日本語でも中国語でもかなり優秀な単語を切るシステムが開発されてきているが、最終的にはやはり「人の手」が頼りである。ここがアルファベット言語との大きな違いである。

また、中国語の場合、Character (字)、Word (単語)、Phrase (句)、Sentence (文) をどう分けるかの問題もかなり厄介である。現在、中国国内でも以下のような「漢語分詞系統」があるが、現代語、近代語、古典語でも変わってくるし、やはり一筋縄ではいかないものであり、中国語学の専門家との「協働」も不可欠になっている。

THULAC <http://thulac.thunlp.org>

HanLP <https://github.com/hankcs/HanLP>

NLPIR <http://ictclas.nlpir.org>

この他、中国語の場合、漢字の字体も問題となる。簡体字・繁体字・異体字をどう処理するかである。近年多くの電子ブックが世に出回っているが、テキストの信頼性に問題のあるものが少なくない。

例えば、「罷／吧」「很／狠」「里／裡／裡／裏」「您／您」など、これらは違いがその成書年代とも関わってくるものであり、それを無視して一つの漢字で表記されたのでは、全く研究には使い物にならないのである。

6. 多くの可能性——東アジア研究の Hub として

ところで、私たちの KU-ORCAS の研究ユニットとその主な内容は以下のとおりであるが、この他、広く学内外からの研究ユニットも募集している。

[ユニット1：東西文化接触とテキスト]

本学所蔵の東西言語接触に関わる資料（辞書・文法書・宣教師報告書等）を中心としたアーカイブ。また、本学所蔵書のほか大英図書館・フランス国立図書館・バチカン図書館・ハーバード大学など海外諸機関の蔵書を相互リンクによって統合したものを構想する。

[ユニット2：東アジアの中の大坂の学統とネットワーク]

本学の学統の源流たる「泊園書院」に関する総合アーカイブを構築する。

また、本学が集中的に所蔵する近世大坂画壇コレクションを中心に国内・海外に散在する大坂画壇作品を含めたデジタルアーカイブを構築する。

[ユニット3：古都・史跡の時空間]

高松塚古墳の発掘に象徴される本学の古代飛鳥・難波津研究が蓄積してきた発掘データ・出土物データ・図面等をデータベース化するとともに、新たに飛鳥時代の墳墓の調査を行い、成果展覧会を開催し、これらを総合したアーカイブを構築する。京都の郊外都市・淀川流域の古文書・古地図・寺社境内絵図を調査・デジタル化する。

なお、こうした研究ユニットは決して「閉じた」ものではなく、私たちのKU-ORCASの最大の特徴は上述の「サイロ問題の打破」を実現すべく、「コンテンツを解き放す」ことを目標にしていることである。

具体的には研究リソースのオープン化、研究グループのオープン化、研究ノウハウのオープン化であるが、全ての東アジア研究者（実は研究者に止まらず、一般市民、学生等々様々なステークホルダーを念頭に置いている）に利用可能な研究プラットフォームを提供することである。

今後、様々な機関と連携しながら、東アジア研究オープン・プラットフォームの実現と知識基盤社会に適合した新たな人文知の構築に力を注いでいきたいと考えている。