

人文学における オープンデータの活用

一般財団法人人文情報学研究所 主席研究員
東京大学大学院情報学環 特任准教授
永崎研宣

「オープンデータ」

- 「自由に使える再利用もでき、かつ誰でも再配布できるようなデータのことだ。従うべき決まりは、せいぜい「作者のクレジットを残す」あるいは「同じ条件で配布する」程度である。」

- <http://opendatahandbook.org/ja/what-is-open-data/index.html>

- CC BY あるいは CC BY SA ?

人文学におけるデータの活用

- 多くは非オープンデータ。
 - 著作権保護されている資料のため個人で入力or OCR or 何らかの方法で入手した(主にテキスト)データ
 - データ自体の配布・再利用不能
 - ここから出てきた成果の検証も不可能
 - 一定の条件下で利用を許可されたデータ
 - 学術利用のみ
 - 再配布不可
 - 一定の条件下で再配布も許可されたデータ
 - 学術利用のみ
 - 商用利用不可

人文学向けのオープンデータ？

- 米国人文学科学基金(NEH)の過去の助成金リスト
 - Data.govからオープンデータとして公開されている
 - 助成金のデータから米国人文学の研究動向が確認できる
 - ……というのはともかくとしまして。
- 欧米では徐々にオープンデータの資料が広まってきている。

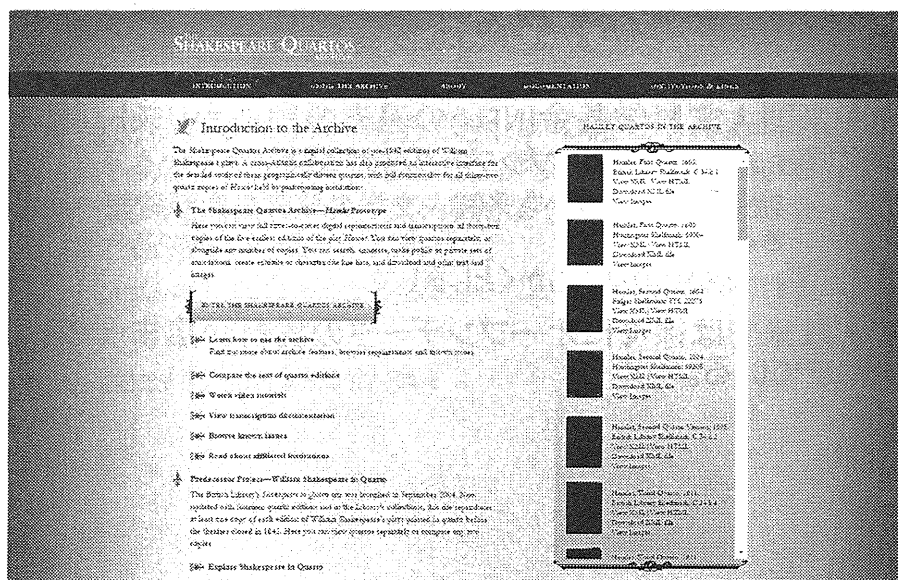
欧米の人文学向けオープンデータ

• プラットフォームレベル

- Europeana
- Hathitrust / Google Books / Internet Archive
- 統合検索エンジンとしてのDPLA
- Gallica @BnF

• 個別の資料レベル

- シェイクスピア・アーカイブ
 - 各地の図書館の連携により30以上の版本の構造化テキスト作成+画像公開
 - ⇒専門家による深い作り込み/ CC BY-NC
- 各種クラウドソーシング翻刻プロジェクト



TEI/XMLによる詳細なマークアップ

[illegible]

内容に基づく様々な構造化

```

5774 </p>
5775 <!-->
5776 <!-->
5777 <!-->
5778 <!-->
5779 <!-->
5780 <!-->
5781 <!-->
5782 <!-->
5783 <!-->
5784 <!-->
5785 <!-->
5786 <!-->
5787 <!-->
5788 <!-->
5789 <!-->
5790 <!-->
5791 <!-->
5792 <!-->
5793 <!-->
5794 <!-->

```

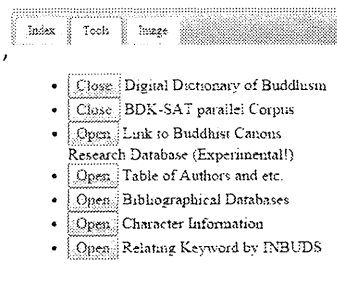
- 「ある登場人物の発言すべて」「Aに続くBの発言のみをすべて」などの様々な情報をTEI/XMLの構造を利用して簡単に抽出可能。

一方、我が国では...

- オープンデータと呼べるものはきわめて少ない
 - 「東寺百合文書」は今のところ例外的
 - CC BYで公開
 - 青空文庫は素晴らしいオープンデータだが研究者向けとまでは言えない。
 - 日本語版Wikipedia は成果の発信に有効
- それでも可能性は様々に垣間見える
 - 「翻デジ」の可能性
 - 仏典のデータベースを参考に。
 - 一部、オープンデータを活用している。

例：大蔵経テキストデータベース

- 様々な情報資源と連携することで利便性を提供している。
 - 電子仏教辞典
 - 漢文・英訳対照コーパス(仏教伝道協会)
 - 漢字データベース関連
 - CHISE(京都大学), HNG(北海道大学等), Unihan, HMS...
 - 書誌データベース
 - INBUDS(日本印度学仏教学会), SARDS3(ドイツ・ハレ大学), CiNii(NII)
 - 目録データベース
 - BCRD(コロンビア大学)
 - 文書画像データベース
 - 近デジ、国デコ(NDL)、Gallica(フランス)、BL(イギリス)、国文学研究資料館DB、その他各大学画像DB



2. 最長一致で自動的に分割して個々の意味を辞書から検索

1. テキストをドラッグすると...

T0262_09.0001b27. 之。金河難命。道殊半滿之科。豈非教被棄
T0262_09.0001b28. 時。無足最其高會。是知五
T0262_09.0001b29. 億之億。五百授記。俱榮
T0262_09.0001c01. 光現瑞。開發諸之教源。出
T0262_09.0001c02. 宏略。朽宅通入大之文軌
T0262_09.0001c03. 隆。解殊明理性之常存。
T0262_09.0001c04. 5 鈎輦宛然。喻障惟遠。
T0262_09.0001c05. 溺之沈流。一極悲心。拯也
T0262_09.0001c06. 唐六百餘載。總歷詳經四千餘部。支何益乎。
T0262_09.0001c07. 無出此經。將非機教相和。並皆勝之遠塵。
T0262_09.0001c08. 聞而深敬。俱感王之餘勸。觀於經首。序而綜
T0262_09.0001c09. 之。庶得早淨六根。仰慈尊之嘉會。速成四
T0262_09.0001c10. 永貽諸後。云
T0262_09.0001c11. 永貽諸後。云
T0262_09.0001c12. 永貽諸後。云
T0262_09.0001c13. 永貽諸後。云
T0262_09.0001c14. 永貽諸後。云
T0262_09.0001c15. 永貽諸後。云
T0262_09.0001c16. 永貽諸後。云
T0262_09.0001c17. 永貽諸後。云
T0262_09.0001c18. 永貽諸後。云
T0262_09.0001c19. 永貽諸後。云
T0262_09.0001c20. 永貽諸後。云
T0262_09.0001c21. 永貽諸後。云
T0262_09.0001c22. 永貽諸後。云
T0262_09.0001c23. 永貽諸後。云
T0262_09.0001c24. 永貽諸後。云
T0262_09.0001c25. 永貽諸後。云
T0262_09.0001c26. 永貽諸後。云
T0262_09.0001c27. 永貽諸後。云
T0262_09.0001c28. 永貽諸後。云
T0262_09.0001c29. 永貽諸後。云
T0262_09.0001c30. 永貽諸後。云
T0262_09.0001c31. 永貽諸後。云
T0262_09.0002a01. 千人俱。1 羅羅羅羅。2 那那羅羅比丘尼。亦
T0262_09.0002a02. 與眷屬俱。菩薩摩訶薩八萬人。皆於阿彌
T0262_09.0002a03. 多羅三藐三菩提不退轉。皆得。3 陀羅尼畢
T0262_09.0002a04. 4 辯才。轉不退轉法輪。供養無量百千諸

Digital Dictionary of Buddhism

電子佛敎辭典

パスワードがない場合は「guest」でログインしてください。
Users who do not have a password can log in with the user ID "guest".

検索語: 如是我聞。一時佛住王舍城耆闍崛
山中。與大比丘衆萬二千人俱。

如是我聞 thus have I heard (rúshì wǒ wén)
一時 one time, at the same time (yí shí)
佛住 a buddha-abode (fózhù)
王舍城 Rājagṛha (Wángshè chéng)
耆闍崛山 Gṛdhraṁṭha (Qíshèjué shān)
衆 aśuddha (zhòng)
與大比丘衆 with a great assembly of monks
(yǔ dà bīqū zhòng)
萬 ten thousand (wàn)
二千 two thousand (èrqiān)
人 human being (rén)
俱 together with (jù)

如是我聞 thus have I heard (rúshì wǒ wén)
一時 one time, at the same time (yí shí)
佛住 a buddha-abode (fózhù)
王舍城 Rājagṛha (Wángshè chéng)
耆闍崛山 Gṛdhraṁṭha (Qíshèjué shān)
衆 aśuddha (zhòng)
與大比丘衆 with a great assembly of monks
(yǔ dà bīqū zhòng)
萬 ten thousand (wàn)
二千 two thousand (èrqiān)
人 human being (rén)
俱 together with (jù)

一時佛住王舍城耆闍崛山中。與大比丘衆萬二千人俱。皆是
阿羅漢。結集已盡。無復煩惱。速得已利盡諸有結。心得自
在。(妙法蓮華經)

1. テキスト中の単語をドラッグすると...

2. 選択した項目の詳細が表示される。

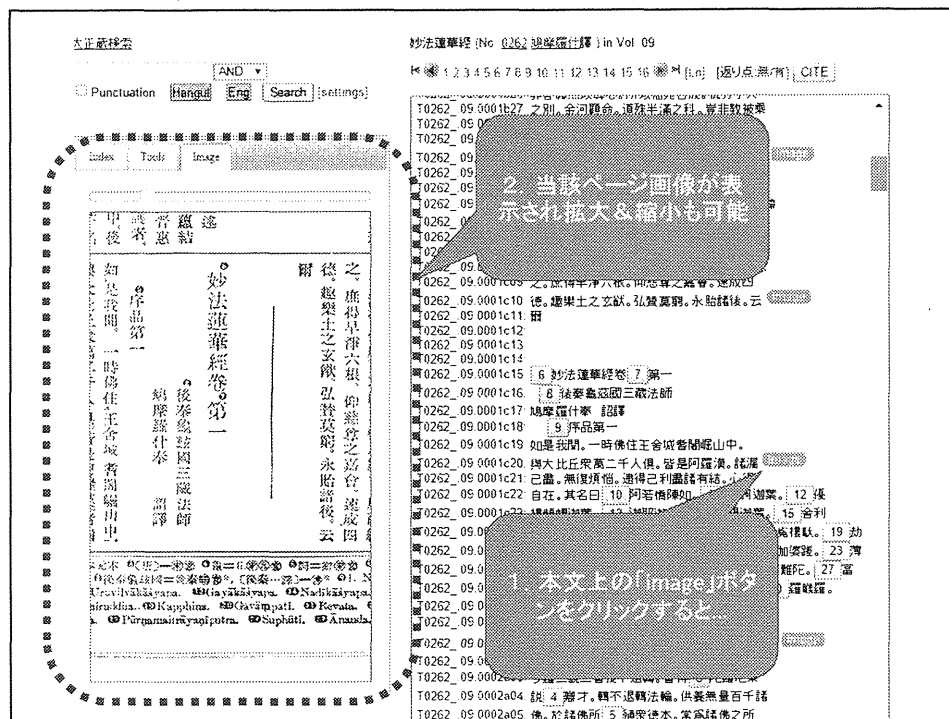
T0262_09.0001b27. 之。金河難命。道殊半滿之科。豈非教被棄
T0262_09.0001b28. 時。無足最其高會。是知五
T0262_09.0001b29. 億之億。五百授記。俱榮
T0262_09.0001c01. 光現瑞。開發諸之教源。出
T0262_09.0001c02. 宏略。朽宅通入大之文軌
T0262_09.0001c03. 隆。解殊明理性之常存。
T0262_09.0001c04. 5 鈎輦宛然。喻障惟遠。
T0262_09.0001c05. 溺之沈流。一極悲心。拯也
T0262_09.0001c06. 唐六百餘載。總歷詳經四千餘部。支何益乎。
T0262_09.0001c07. 無出此經。將非機教相和。並皆勝之遠塵。
T0262_09.0001c08. 聞而深敬。俱感王之餘勸。觀於經首。序而綜
T0262_09.0001c09. 之。庶得早淨六根。仰慈尊之嘉會。速成四
T0262_09.0001c10. 永貽諸後。云
T0262_09.0001c11. 永貽諸後。云
T0262_09.0001c12. 永貽諸後。云
T0262_09.0001c13. 永貽諸後。云
T0262_09.0001c14. 永貽諸後。云
T0262_09.0001c15. 永貽諸後。云
T0262_09.0001c16. 永貽諸後。云
T0262_09.0001c17. 永貽諸後。云
T0262_09.0001c18. 永貽諸後。云
T0262_09.0001c19. 永貽諸後。云
T0262_09.0001c20. 永貽諸後。云
T0262_09.0001c21. 永貽諸後。云
T0262_09.0001c22. 永貽諸後。云
T0262_09.0001c23. 永貽諸後。云
T0262_09.0001c24. 永貽諸後。云
T0262_09.0001c25. 永貽諸後。云
T0262_09.0001c26. 永貽諸後。云
T0262_09.0001c27. 永貽諸後。云
T0262_09.0001c28. 永貽諸後。云
T0262_09.0001c29. 永貽諸後。云
T0262_09.0001c30. 永貽諸後。云
T0262_09.0001c31. 永貽諸後。云
T0262_09.0002a01. 千人俱。1 羅羅羅羅。2 那那羅羅比丘尼。亦
T0262_09.0002a02. 與眷屬俱。菩薩摩訶薩八萬人。皆於阿彌
T0262_09.0002a03. 多羅三藐三菩提不退轉。皆得。3 陀羅尼畢
T0262_09.0002a04. 4 辯才。轉不退轉法輪。供養無量百千諸

2. 対訳コーパスから検索結果を引き出す

BDK-SAT parallel corpus

Search results in BDK-SAT corpus:
If any Chinese text below is clicked, the entire text is displayed
on the center window.
All texts in this corpus are downloadable at BDK Web site in
PDF format.

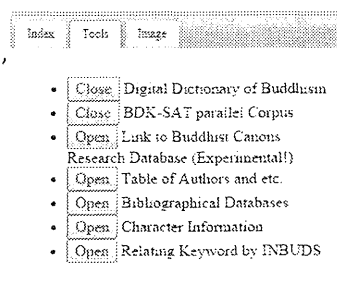
1. 諸比丘等。於法華經。若坐
2. 諸比丘等。於法華經。若坐
3. 諸比丘等。於法華經。若坐
4. 諸比丘等。於法華經。若坐
5. 諸比丘等。於法華經。若坐
6. 諸比丘等。於法華經。若坐
7. 諸比丘等。於法華經。若坐
8. 諸比丘等。於法華經。若坐
9. 諸比丘等。於法華經。若坐
10. 諸比丘等。於法華經。若坐
11. 諸比丘等。於法華經。若坐
12. 諸比丘等。於法華經。若坐



再掲：大蔵經テキストデータベース

- 様々な情報資源「オープンデータ」と呼べそうなものは...？

- 電子仏教辞典
- 漢文・英訳対照コーパス(仏教伝道協会)
- 漢字データベース関連
 - CHISE(京都大学), HNG(北海道大学等), Unihan, HMS...
- 書誌データベース
 - INBUDS(日本印度学仏教学会), SARDS3(ドイツ・ハレ大学), CiNii(NII)
- 目録データベース
 - BCRD(コロンビア大学)
- 文書画像データベース
 - 近デジ、国デコ(NDL)、Gallica(フランス)、BL(イギリス)、国文学研究資料館DB、その他各大学画像DB



「オープンアクセス」では？

- PDFによる論文公開が分野としては徐々に広がってきている
 - 仏教学では:
 - 日本印度学仏教学会
 - Journal@rticle + NII-ELS ⇒ JSTAGE or JSTAGE Lite
 - それ以外にも:
 - 主にNII-ELS経由
 - 日本語学会:
 - CiNiiにてOA
 - ⇒ NII-ELSの多大な貢献
- Open Library of Humanitiesの取り組み
- OpenEdition.orgの取り組み

データ作成者への評価

- 「校訂テキスト」「目録」「事典」等とのアナロジーはどうか
 - 同分野の研究者コミュニティからの評価は、冬の時代を迎えつつある現在、どの程度有効か。
 - 他の研究分野からも評価される枠組みが必要
 - 人件費/謝金はどのようにカバーし得るか
 - プロの編集は必要か・必要ならどうするか
- データジャーナル的なものが評価指標たり得るか
 - データの質についてはどう評価するか
 - プログラムの動作確認等と異なり、人文学向けデータの内容の確認の自動化はかなり困難。

終わりに

- 人文学向け資料では「オープンデータ」はまだまだ少ない。
 - 特に日本では弱いような印象があるが引き続き調査したい。
 - エコシステムとして機能するに至っていない？
 - オープンデータでなくとも研究用途ではそれなりに有益。
 - しかし、商用利用を禁ずることで広がりを阻害している面もあるかもしれない。
- 我が国での人文学向け資料における「オープンデータ」を活用したソリューションの必要性