

大学・研究機関による データキュレーションの実践に 向けて

国立情報学研究所 南山 泰之

琉球大学 研究データ管理セミナー 2024年2月15日

(於) 琉球大学附属図書館2階ラーニング・コモンズ



目次

1. はじめに

2. 国内の動向

3. 実践に向けて

4. まとめ



1. はじめに



データキュレーションとは

用語の関係性: Curation > Archiving > Preservation

- **キュレーション**: データが目的に適合し、発見や再利用できるようにするために、データが作成された時点から、データの利用を管理し促進する活動。動的なデータセットの場合、これは目的に合ったデータを維持するための継続的なエンリッチメントまたは更新を意味する。より高度なレベルのキュレーションには、アノテーションや他の出版物とのリンクを維持することも含まれる
- **アーカイビング**: データが適切に選択され、保管され、アクセスできるようにし、セキュリティと真正性を含め、論理的および物理的な整合性が長期間にわたって維持されるようにするキュレーション活動
- **プリザベーション**: データの特定の項目が時間の経過とともに維持され、技術の変化によってもアクセスして理解できるようにするための<u>アーカイビングの活動</u>



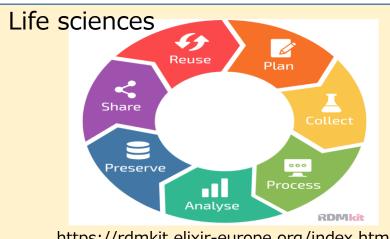
データを解釈可能かつ利用可能にするための一連の処理を **データ**キュレーションと呼ぶ

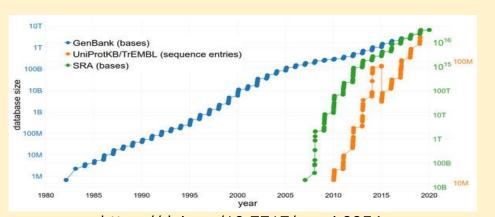
Ref: e-Science Curation Report: Data curation for e-Science in the UK: an audit to establish requirements for future curation and provision prepared for: The JISC Committee for the Support of Research (JCSR). 2003 http://digitalpreservation.gov/news/2004/e-ScienceReportFinal.pdf



各分野におけるデータキュレーションの実践

生命科学、社会科学などの分野ではデータキュレーションの ライフサイクルを定め、効率的なデータキュレーションを実践している

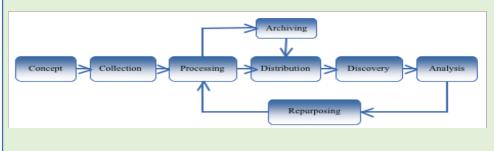


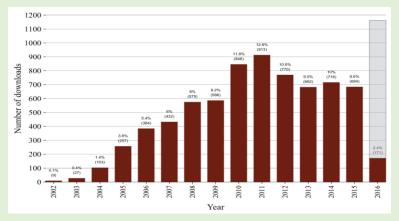


https://rdmkit.elixir-europe.org/index.html

https://doi.org/10.7717/peerj.9954

Social sciences





https://ddialliance.org/Specification/DDI-Lifecycle/

https://doi.org/10.1177/2158244016685136



欧州委員会(EC)によるFAIR原則の費用便益分析

(1) 費用便益分析で対象とした要素と損失値

領域	要素	
研究活動へのインパクト	(a) 費やす時間 (b) 保管コスト (c) ライセンスコスト	_
将来の研究の機会へのインパクト	(d) 研究の撤回率 (e) 助成の重複 (f) 学際性の欠如	1
イノベーションへのインパクト	(g) 潜在的な経済成長	

少なくとも **年間102億ユーロ以上**の損失

※【試算方法】(a)~(e): FAIR原則に従う研究データがないときの非効率性(費やす時間と関連コスト)を評価。研究者が研究データへアクセスするために支払う必要のあるライセンスのコストを見積。FAIR原則が実施されていれば不必要である重複した研究データの保管コストを算出。(f)と(g): 定性的な考察に基づく。



(2) 上記の費用便益分析で対象としていない要素と損失の推定値

要素

研究の質 経済的取引高 研究データの機械可読性



推定**年間160億ユーロ**の損失

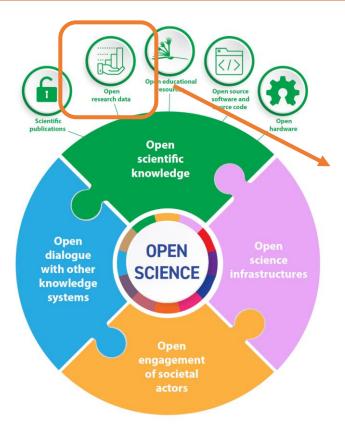


※FAIR原則:研究データの適切な共有・公開を行う もっともよく知られた基準。学術コミュニティにおけ るデータキュレーション実践のための事実上の指針。

FAIR原則に沿ったデータキュレーションが行われない際の損失は**年間262億ユーロ**



ユネスコ「オープンサイエンスに関する勧告」



"Open research data are available in a timely and user-friendly, human- and machine-readable and actionable format, in accordance with principles of good data governance and stewardship, notably the FAIR (Findable, Accessible, Interoperable, and Reusable) principles, supported by regular curation and maintenance"

FAIR原則に基づくデータキュレーションが社会的にも求められつつある

Ref: UNESCO Recommendation on Open Science https://unesdoc.unesco.org/ark:/48223/pf0000379949.locale=en

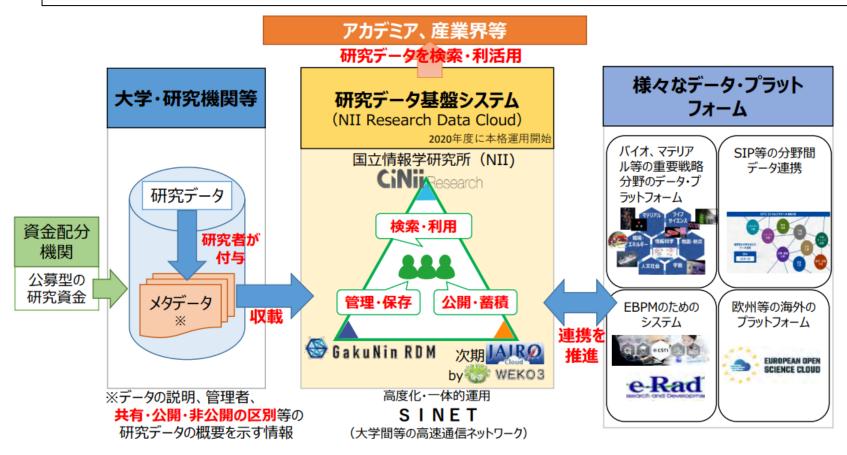


2. 国内の動向



政策的な要請

研究データ基盤システムを中核としたデータ・プラットフォームの構築 →研究データの公開・共有を推進、産学官のユーザがデータを検索可能



Ref: 研究データの管理・利活用に関する現状と課題について

https://www8.cao.go.jp/cstp/gaiyo/yusikisha/20220428/siryo3_1_1.pdf



FAIR原則への言及

「研究DXの推進一特にオープンサイエンス、データ利活用推進の視点からーに 関する審議について | (2022年12月23日)

内閣府からの審議依頼を受けて、日本学術会議に設置したオープンサイエンスを推進するデータ 基盤とその利活用に関する検討委員会、同オープンサイエンス企画分科会及び同オープンサイエン ス企画分科会オープンサイエンス・データ利活用推進小委員会が中心となり審議を行った。

【提案1】研究者が容易に利用可能な研究データプラットフォームの構築

【提案2】データプロフェッショナルの育成と多面的な研究評価の実現

【提案3】モニタリング機構に基づくデータ駆動型研究の不断の改善

【提案4】研究自動化(ARW)に向けた情報技術、計算資源の集約

【提案5】分野を越えた連携を実現する FAIR 原則の追求

【提案6】法制度面でのデータガバナンスの構築

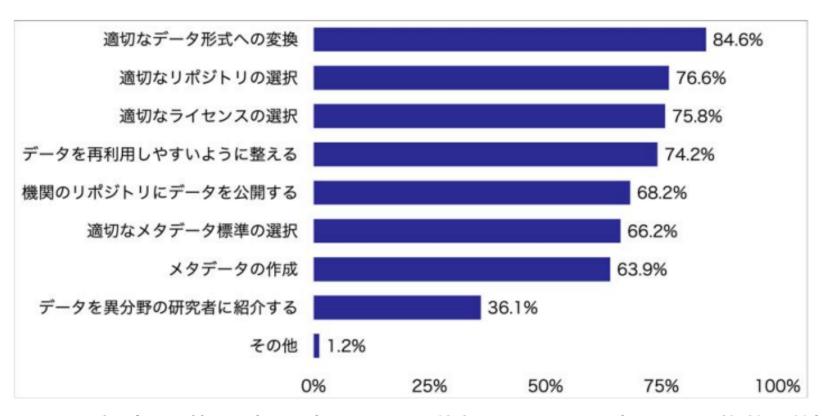


"多様な分野間の連携においては論文だけではなくデータを介することによりARW の更なる加速が期待されるが、いかなるメタデータ付与がデータの利活用に資する かは分野依存性が高い。FAIR 原則に基づき、各分野における積極的な分野間連携 の実践が望まれ、また、その経験値の共有を可能とする取組が推奨される。日本学 術会議や学会等において、FAIR 原則を追求する踏み込んだ議論が期待される。

> Ref: 回答 研究 D X の推進 – 特にオープンサイエンス、データ利活用 推進の視点からーに関する審議について https://www.scj.go.jp/ja/info/kohyo/pdf/kohyo-25-k335.pdf 10



研究者からの要望



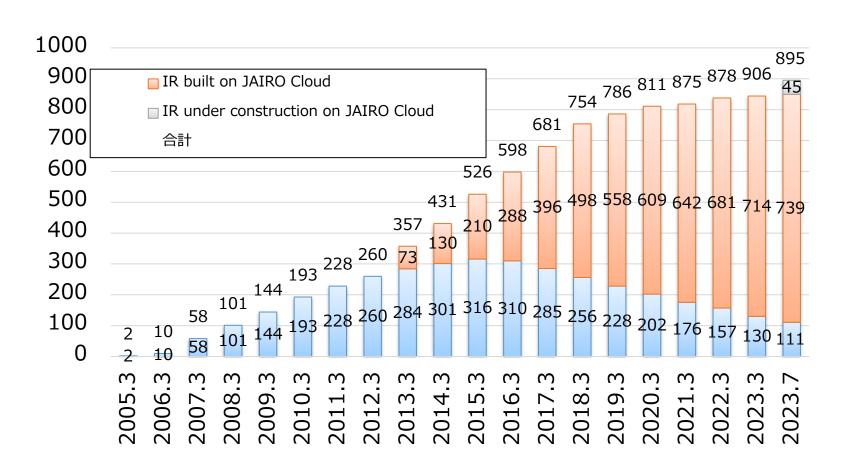
研究データ管理(RDM)に関して依頼したい項目(n=488:複数回答) 図 61

データ利活用に向けた要望が高い

Ref: 池内 有為, 林 和弘. 研究データ公開と論文のオープンアクセスに関する実態調査2020 https://doi.org/10.15108/rm316



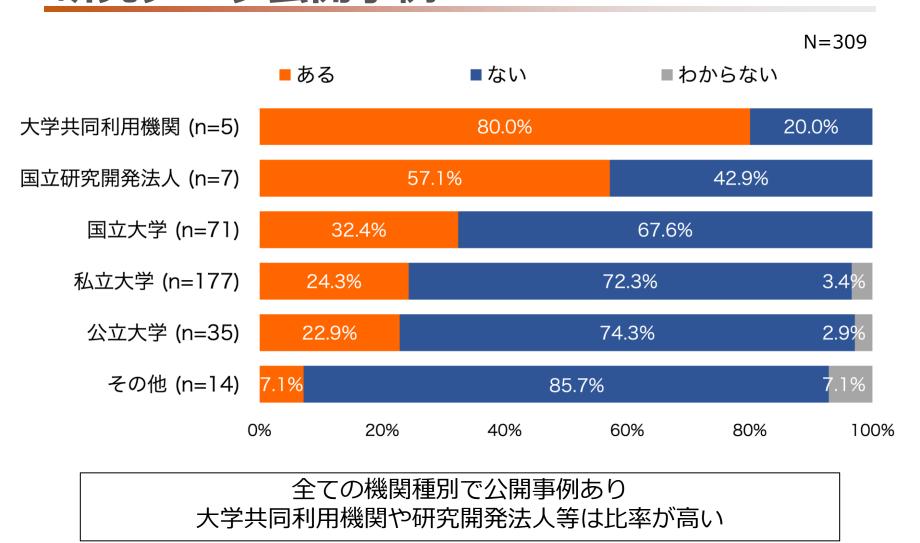
大学・研究機関の準備体制



国内の機関リポジトリ数は895を数える

80%以上がJAIRO Cloudサービスを利用

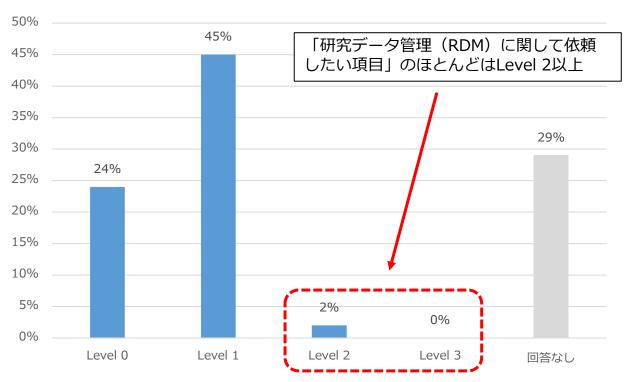
2021~2022年の機関リポジトリによる Parkers and Data Platform 研究データ公開事例



Ref: 「国内機関における研究データ管理の取り組み状況調査2022 | 集計結果より 13



課題:データ公開に向けた作業レベル



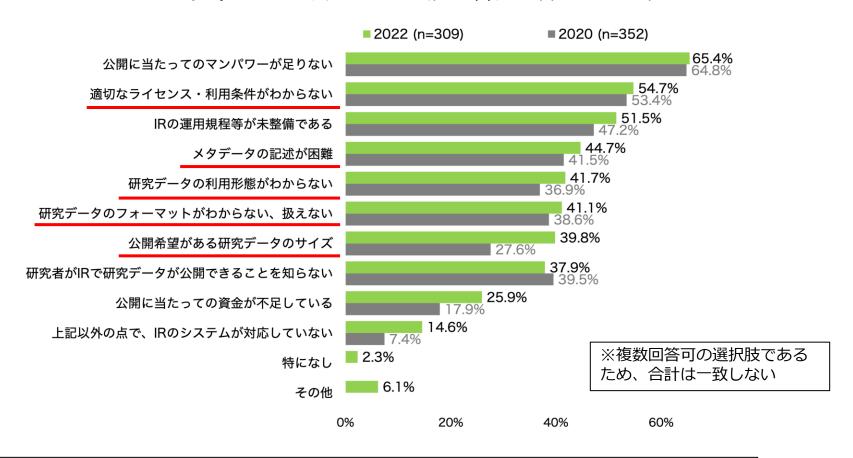
レベル	見出し	説明(抜粋)
Level 0	登録データの流通	登録されたままで公開する
Level 1	基礎的な標準に準拠	簡単なチェック、基礎的なメタデータや文書の追加を行う
Level 2	論理・技術面での検証	フォーマット変換、文書やメタデータの強化などを実施する
Level 3	概念的な理解と再利用	Level 2に加えて、データレベルでの編集を行う

Ref: 「機関リポジトリ/データリポジトリの運用実態に関するアンケート調査報告書」 https://jpcoar.repo.nii.ac.jp/records/2000154



何が課題や障壁となっているのか

機関リポジトリ(IR)による研究データ公開の課題や障壁となり得ること



研究分野によって作法が異なる項目が多く含まれている

Ref: 「国内機関における研究データ管理の取り組み状況調査2022」集計結果より 15

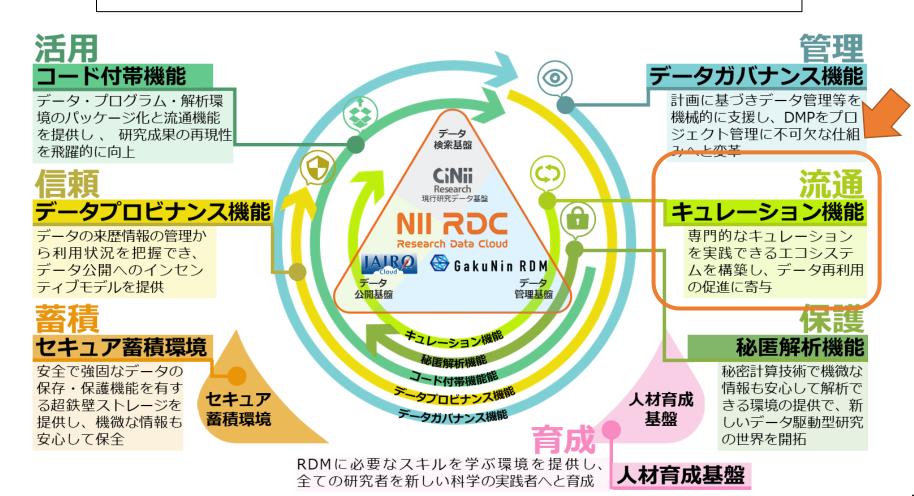


3. 実践に向けて



データキュレーション支援機能

汎用的なプロセスをアプリケーションとして実装し、 大学・研究機関へデータキュレーションサービスを提供



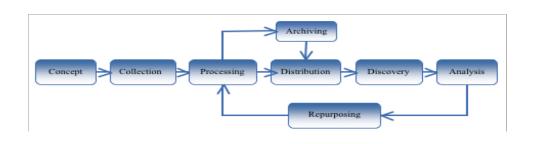


課題:各分野における捉え方の違い

Life sciences



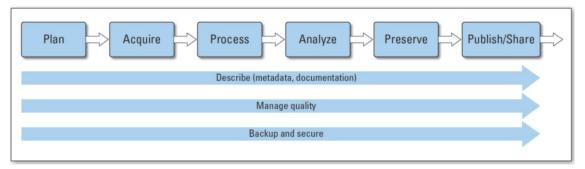
Social sciences



https://rdmkit.elixir-europe.org/index.html

https://ddialliance.org/Specification/DDI-Lifecycle/

Earth science



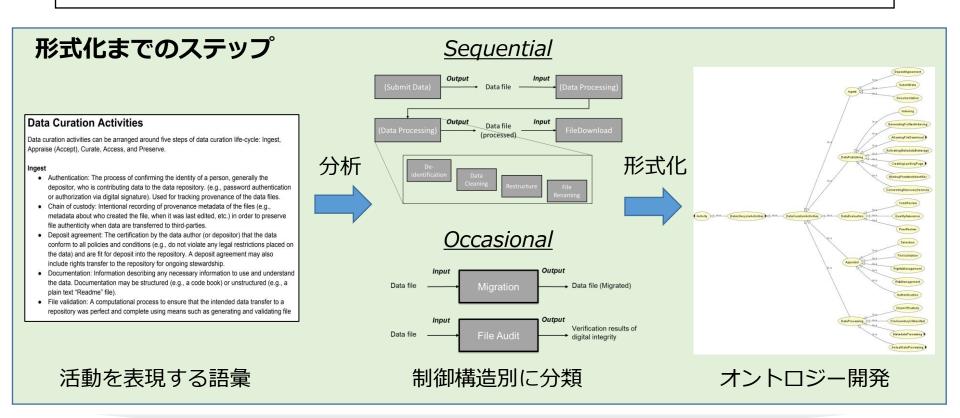
https://doi.org/10.3133/ofr20131265

各分野で定義されるライフサイクルは、分野内での再利用に最適化されている



データキュレーション活動の形式化

複数分野におけるデータキュレーション活動の実態を調査し、 活動に含まれるプロセスを形式化



分野横断的なデータキュレーションプロセスを抽出



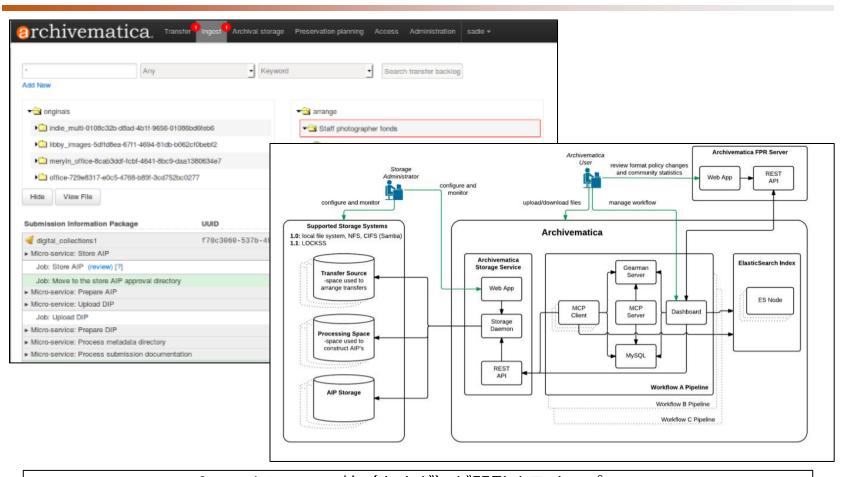
NII RDCとの対応付け

			GakuNin	JAIRO		
カテゴリ	活動名称	対応する機能要件(短縮版)	RDM	Cloud		
		あらかじめ定められたパッケージ規格に合わせてフォルダ構				
	配置・記述	造を変換できること	×	X		
		アナログ形式のデータファイルをデジタル形式に転換できる				
	コンバージョン	こと	×	X		
		提出されたデータファイルに含まれる欠陥やエラーを検出で				
	データクリーニング	きること	×	X		
	データの再構造化	提出されたデータファイルの構造を整理できること	×	X		
		提出されたデータファイルをもとにしたグラフ・図・表を作				
	データ可視化	成できること	×	0		
		デークファフリ に今まわてわ、シー・ゴ+ン1桂却(/田 L 桂却				
	ロカル	データファイルに含まれるセンシティブな情報(個人情報、	V	\ <u>\</u>		
DataProc essing	匿名化	要保護情報など)を書き換えたり削除したりできること	X	×		
	コーノリフューフット亦悔	データをオープンで非独占的なファイル形式に変換できるこ	V	\ <u>/</u>		
	ファイルフォーマット変換		X	X		
	リネーム	データのファイル名を変更できること	0	X		
	相互運用性	分野別の標準を用いてデータをフォーマットできること	X	X		
	=	一連の処理を通じて、データ一式の履歴情報を記録できるこ				
	証拠保全		0	×		
	7 - /11 /2 7 · 0 6 2	データー式を定期的に点検し、データ数、ファイルタイプ				
	ファイルインスペクション	(拡張子)、ファイルサイズ等を把握できること	×	X		
	->-+	データー式を、関連する出版物、論文、プロジェクト等とリ				
	コンテキストの付与	ンクさせられること	X	×		
	メタデータ生成	所定のスキーマに沿ったメタデータを生成できること	0	0		

データ加工に関する活動のカバー率が低い



ソフトウェア候補の選定

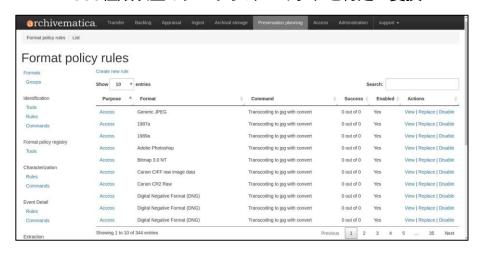


Artefactual Systems社(カナダ)が開発するオープンソース OAIS参照モデルに準拠し、情報パッケージを作成可能な統合ソフトウェアツール群 分野別のユースケースに合わせた複数パイプラインの制御が可能

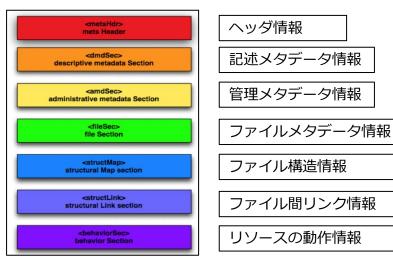


Archivematicaができること

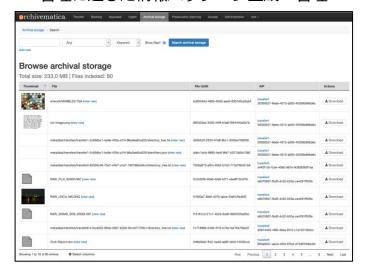
1000種類以上のデータフォーマットを特定・変換



入力ファイルに対してMETSメタデータを自動生成



管理に適した情報パッケージ生成・管理



https://www.loc.gov/standards/mets/METSPrimer.pdf https://www.archivematica.org/en/d

ocs/archivematica-1.15/



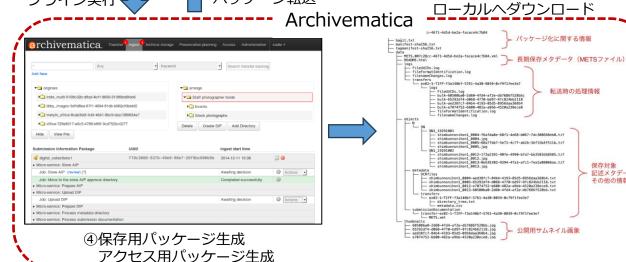
転送時の処理情報

記述メタデータ その他の情報

パイプライン制御UIの開発



- ユーザは制御UIを用いて各種リポ ジトリ(GitHub等)から対象デー 夕を取得できる
- 分野別パイプラインはユーザごと に設定可能
- ユーザ管理はWEKO3と連動して実 施可能



NII RCOS

人文学データのケーススタディ (データエコ事業ユースケース)

- 渋沢栄一記念財団の文書データ
 - 役職員が在職中に執筆した原稿
 - 渋沢栄一に関する論考や、財団の活動内容を含む記録として長期保存したい
- 文書に多く用いられるフォーマット
 - Microsoft Office (Excel, PowerPoint, Word)
 - プレーンテキスト(txt)
- 階層構造を持つ、合計10ファイルを入力した。

```
masahide DA 1
   —metadata
        metadata.csv
   ┗○○執筆原稿
      ├─○○センター保存分
         ─200604_慶応大学講座
              2006course.doc
              慶応大学記念講座開講のお知らせ.txt
              慶応講座、栄一,4.11.06.doc
              添付書類の訂正.txt
         └-201107 新渡戸国際塾
               20110709○○理事長新渡戸国際塾講義 01.pptx
               readme02.txt
               新渡戸P P.docx
               新渡戸11.docx
      一管理用リスト
            理事長講演と執筆リスト.xlsx
            理事長講演と執筆リストrk(更新用).xlsx
```



テスト結果

く作成されたAIPの中身>

```
AIP作成後のMETS (PREMIS, DCなどのメタデータを格納)
data
    METS.5fc9f5c4-6f46-4cdf-b00f-7c7b35dff029.xml
    RFADMF.html
                                                  処理ログ
  ⊬logs
       fileFormatIdentification.log
       filenameChanges.log
       FileUUIDs.log

—transfers

        2023-10-26 test4-45a464e7-521d-4643-8ee4-613445f4237d
             └-logs
                 -bulk-2281bbc8-a07e-49e4-8d15-21e7fbeebb0f
                ─bulk-266f11eb-752e-4212-a963-40c9d741cb42
                -bulk-27380689-38f8-46b4-9cfd-18cff03b87b3
                ─bulk-2ee70f30-126a-48ce-bc99-9a73c6c10f0f
                -bulk-3b459868-46e1-4f38-8291-8dbb84bbe5a9
                ─bulk-53a03c61-81f0-41f4-88de-317f0572c051
                bulk-8718732f-1c39-4027-a555-3aac307de1ad
                ─bulk-a1d321eb-601b-45de-aa79-2e8dd3e4c486
                ─bulk-bcb133ca-345b-457b-ab84-48106c7d1aba
                bulk-be67276f-3497-4852-b516-b905a4d58b0d
```

```
受け入れたデータの初期状態と入力したメタデータ
└─objects
   H metadata

—transfers

          2023-10-26 test4-45a464e7-521d-4643-8ee4-613445f4237d
                                                                   長期保存用
                  directory tree.txt
                                                                   データ
                  metadata.csv
   ├─Shibui_Ze_Ya_Ying_Zhi_Bi_Yuan_Gao_
      ├─Guan_Li_Yong_risuto
             Li Shi Chang Jiang Yan toZhi Bi risuto.xlsx
             Li_Shi_Chang_Jiang_Yan_toZhi_Bi_risutork(Geng_Xin_Yong_).xlsx
      └─Qing_Bao_Zi_Yuan_sentaBao_Cun_Fen_
          ├─200604 Qing Ying Da Xue Jiang Zuo
                 2006course.doc
                 Qing Ying Da Xue Ji Nian Jiang Zuo Kai Jiang nooZhi rase.txt
                 Qing_Ying_Jiang_Zuo___Rong_Yi__4.11.06.doc
                 Tian Fu Shu Lei noDing Zheng .txt
          └─201107_Xin_Du_Hu_Guo_Ji_Shu_
                 20110709Shibui Ze Li Shi Chang Xin Du Hu Guo Ji Shu Jiang Yi 01.pptx
                 readme02.txt
                 Xin_Du_Hu_11.docx
                 Xin_Du_Hu_P_P.docx
     -submissionDocumentation
                                                                        AIP作成前の
       transfer-2023-10-26 test4-45a464e7-521d-4643-8ee4-613445f4237d
                                                                       METS
              METS.xml
```

全てのログを残した形で階層構造を持つデータのAIP作成に成功 今後、大学図書館におけるユースケース開発に着手予定



4. まとめ

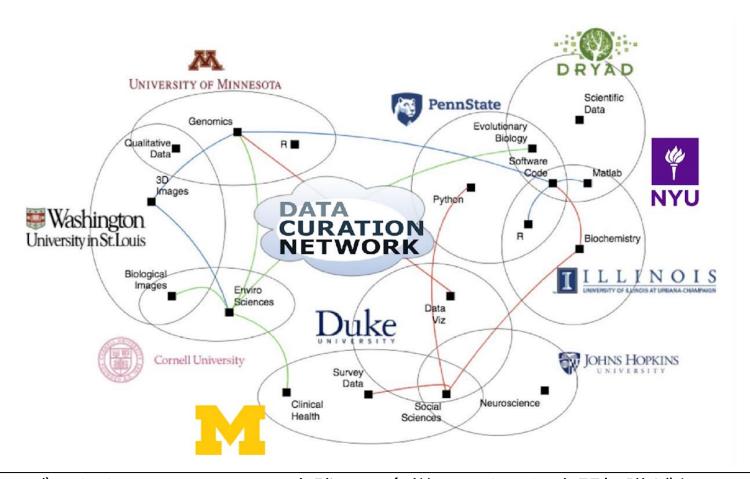


本日の振り返り

- 1. はじめに
 - ✓ データキュレーションとは
 - ✓ 国際的な期待の高まり
- 2. 国内の動向
 - ✓ 政策的な要請/研究者からの期待
 - ✓ 大学側の準備体制
 - ✓ 課題や障壁
- 3. 実践に向けて
 - ✓ データキュレーション支援機能の開発状況紹介
 - ✓ 人文学分野でのケーススタディ
- 4. まとめ
 - 大学図書館との連携を模索中です



データ専門職によるネットワーク構築

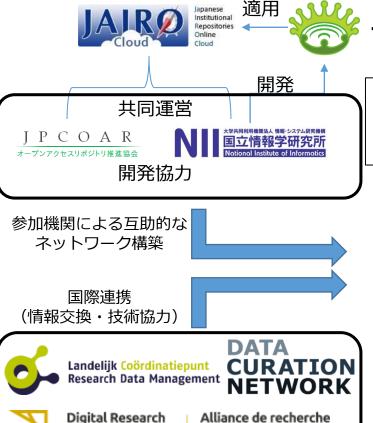


データキュレーションの実践には多様なスキルや専門知識が必要 効率的にデータ整備を行うため、専門職ネットワークの構築が進んでいる



将来構想

どの機関でも専門的なデータキュレーションサービスを 受けられる体制を構築し、データ再利用の促進に寄与



numérique du Canada

Alliance of Canada

データキュレーション 機能

- 大学図書館員、研究支援職員、研究者(データキュレーター)を含む 人的ネットワークを構成
- 単一の機関でカバーしきれない多様な専門分野のデータキュレーションを、大学共同利用機関を含む複数機関間でオンデマンドに支援

