

## GakuNin RDMデータ活用セミナー：これからの研究 データ管理を探る

下山，武司  
国立情報学研究所オープンサイエンス基盤研究センター

藤原，一毅  
国立情報学研究所オープンサイエンス基盤研究センター

清水，敏之  
九州大学データ駆動イノベーション推進本部研究データ管理支援部門

<https://doi.org/10.15017/7159632>

---

出版情報：2023-11-16. National Institute of Informatics  
バージョン：  
権利関係：Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International

# GakuNin RDM データ解析機能の使い方

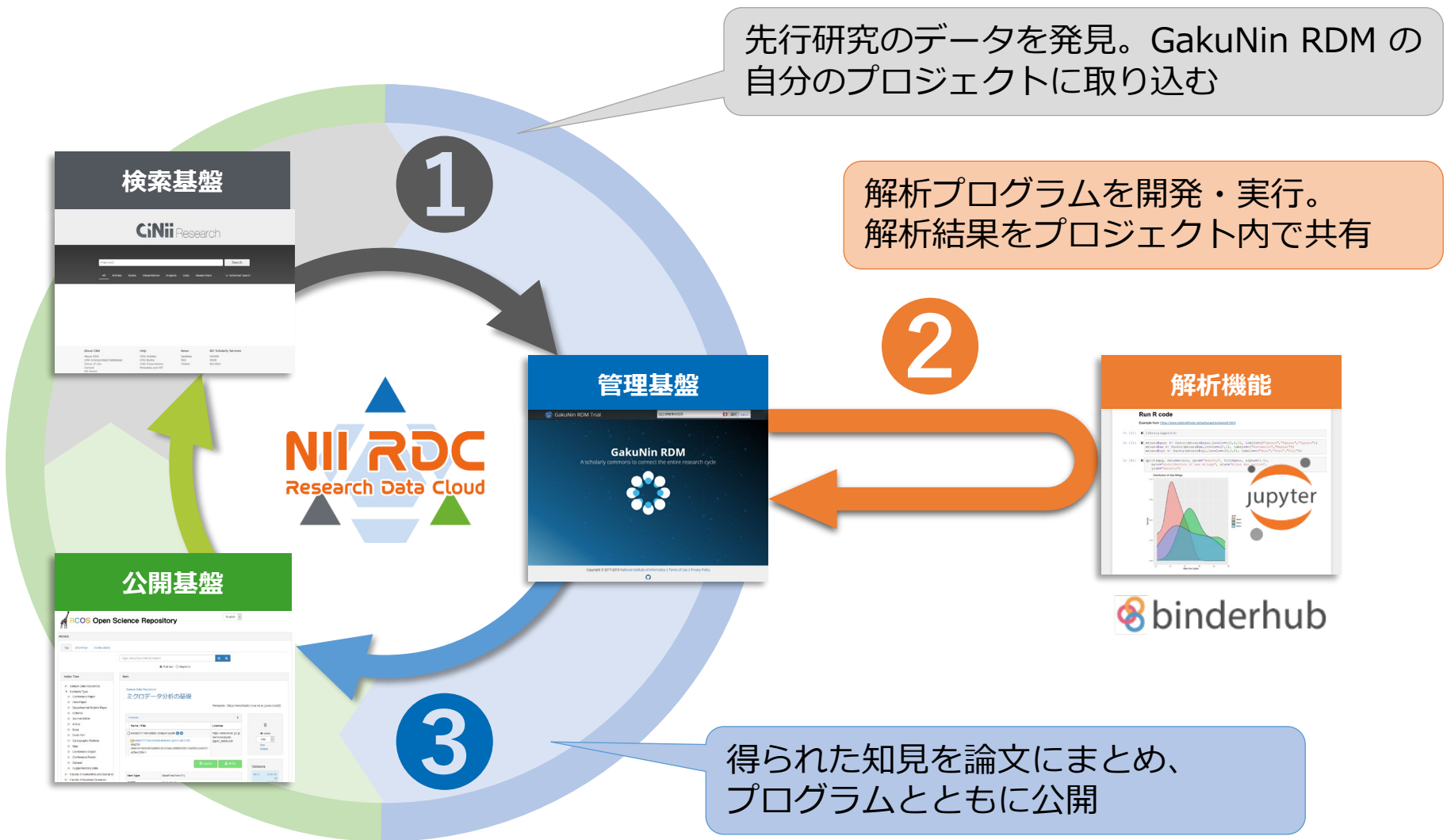
藤原一毅

国立情報学研究所オープンサイエンス基盤研究センター

2023/11/16

九州大学 GakuNin RDM データ活用セミナー

# データとコードが循環する世界



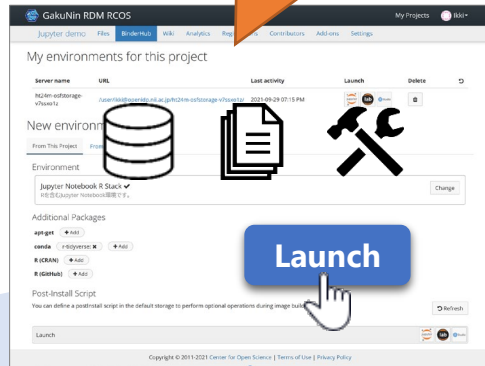
先行研究のデータを発見。GakuNin RDM の自分のプロジェクトに取り込む

解析プログラムを開発・実行。解析結果をプロジェクト内で共有

得られた知見を論文にまとめ、プログラムとともに公開

# GakuNin RDM データ解析機能

①環境定義・共有



GakuNin RDM

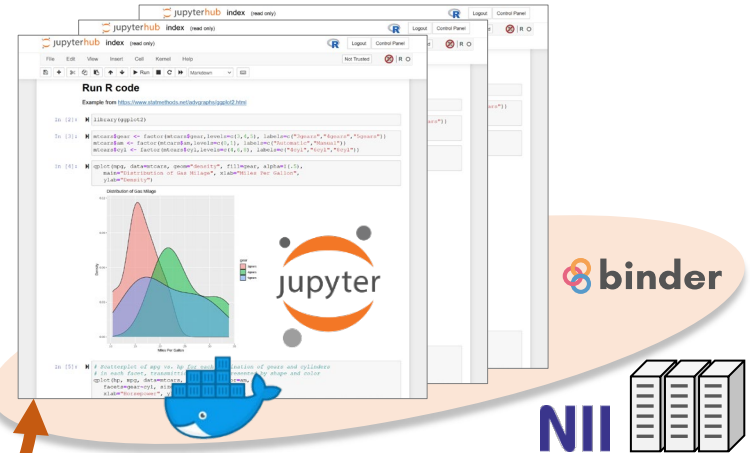
標準ストレージ

機関ストレージ

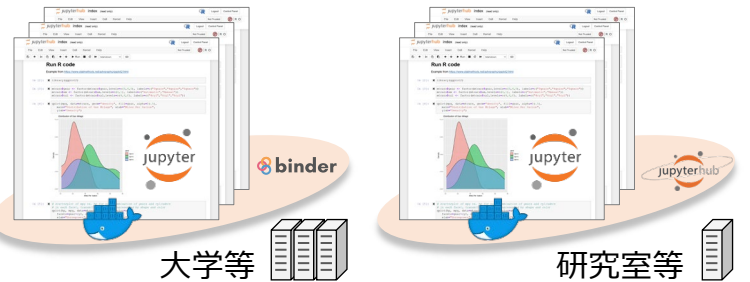
②取り込み

③書き戻し

④読み書き



NII



大学等

研究室等

- JupyterHub がインストールされた計算機と連携し、データ解析環境をGakuNin RDMから1クリックで構築
- NII所有の計算機のほか、クラウド上のVMなど外部計算機とも連携可能

# デモ (1/4)

ブラウザのアドレスバーには `https://rcos.rdm.nii.ac.jp/nc6pd/files/` が表示されています。ページタイトルは「GakuNin RDM RCOS | Reanalysis」。ナビゲーションメニューには「Reanalysis Example」「ファイル」「Wiki」「解析」「メンバー」「アドオン」「設定」「証跡管理」があります。ユーザー名は「Ikki@OpenIdP」です。

ストレージプロバイダーをクリックするか、ドラッグ&ドロップしてファイルをアップロードします

名前	サイズ	バージョン	ダウンロ...	最終更新日時
Reanalysis Example				
- NII Storage				
analyses.ipynb	3.2 kB	1	0	2022-05-27 10:34 AM
supplementary materials.xlsx	150.7 kB	1	0	2022-05-27 10:34 AM

「analyses.ipynb」は「解析プログラム (Jupyter Notebook)」とラベルされています。  
「supplementary materials.xlsx」は「データファイル」とラベルされています。

# デモ (2/4)



## 新しい解析環境

① 解析環境を構成

### 基本イメージ

Python 3.9 + R 4.1.3 ✓

Jupyter Notebook, JupyterLab, RStudio, Shinyが使えます。

変更

### 追加パッケージ

apt-get fonts-noto-cjk: ✕ + 追加

conda seaborn: ✕ openpyxl: ✕ + 追加

pip + 追加

R (MRAN) + 追加

### 自動実行スクリプト

```
#!/bin/bash
set -x
...
```

保存

② 計算機を選択して起動!

### 環境作成

このプロジェクトのデフォルトストレージの内容がコピーされます。

新しい解析環境を作成: <https://binder.cs.rcos.nii.ac.jp>

# デモ (3/4)

**③** ファイルがGakuNin RDMからコピーされている

**②** 書き戻しボタン

**④** ファイルを読み込んで解析

**⑤** 解析結果を ~/result/ に保存

```
[1]: df = pd.read_excel("supplementary_materials.xlsx", index_col=0)
df = df.rename(columns={'day(1=3/31, 2=4/30, 3=5/31, 4=6/10)': 'day'})
df['day'] = df['day'].map({'1': '3/31', 2: '4/30', 3: '5/31', 4: '6/10'})
df

[3]: top10 = pd.pivot_table(df, index='country').nlargest(10, 'Infections')
df1 = pd.pivot_table(df[df['country'].isin(top10)], index='country', columns='day')
ax = df1['Infections'].plot(xlabel='調査日', ylabel='百万人あたり感染者数')
# ax.legend(loc='center left', bbox_to_anchor=(1.0, 0.5))
plt.savefig("result/graph1.png")
```

country	3/31	4/30	5/31	6/10
Belgium	~1000	~2000	~5000	~10000
Chile	~1000	~2000	~5000	~10000
Ireland	~1000	~2000	~5000	~10000
Kuwait	~1000	~2000	~5000	~10000
Luxembourg	~1000	~2000	~5000	~10000
Peru	~1000	~2000	~5000	~10000
Qatar	~1000	~2000	~5000	~10000
Singapore	~1000	~2000	~5000	~10000
Spain	~1000	~2000	~5000	~10000
United States of America	~1000	~2000	~5000	~10000

```
[4]: df2 = pd.read_excel("supplementary_materials.xlsx", sheet_name=4,
index_col=0, header=5, skipfooter=3,
usecols=[1,2,3,5,6,8,9,11,12], skiprows=[6])
df2 = df2.dropna()
df2 = df2.set_axis(['CF1', 'SE1', 'CF2', 'SE2', 'CF3', 'SE3', 'CF4', 'SE4'], axis=1)
df2['3/31'] = df2['CF1'] / df2['SE1']
df2['4/30'] = df2['CF2'] / df2['SE2']
df2['5/31'] = df2['CF3'] / df2['SE3']
df2['6/10'] = df2['CF4'] / df2['SE4']
```

# デモ (4/4)

The screenshot displays the GakuNin RDM RCOS web interface. At the top, the browser address bar shows the URL <https://rcos.rdm.nii.ac.jp/yftz/>. The page header includes the GakuNin RDM RCOS logo, the text "Reanalysis of COVID-19 Infect...", and navigation links for "ファイル", "Wiki", "解析", "メンバー", "アドオン", "設定", and "証跡管理". The user's name "Ikki" is visible in the top right.

The main content area shows a file named "graph1.png (バージョン: 1)". Below the filename are several action buttons: "チェックアウト", "タイムスタンプを打つ", "削除", "ダウンロード", and "プレビュー". A "バージョン管理" button is located below these.

On the left side, a file explorer sidebar is visible, showing a tree view of files and folders. The file "graph1.png" is highlighted in blue and enclosed in a red rectangular box. An orange callout bubble points to this box with the text: "⑦ 解析結果がGakuNin RDMに書き戻される".

The main content area also features a line graph titled "country". The y-axis is labeled "百万人あたり感染数" (Infection rate per 100,000 people) and ranges from 0 to 30,000. The x-axis is labeled "調査日" (Survey date) and shows dates from 3/31 to 6/10. The graph plots data for several countries: Belgium, Chile, Ireland, Kuwait, Luxembourg, Peru, Qatar, Singapore, Spain, and United States of America. Qatar shows the highest and most rapidly increasing infection rate, reaching nearly 30,000 by 6/10. Other countries like the United States of America and Spain also show significant increases, while others remain relatively low.



# こんな用途に使えます

## 研究

- ご自身の研究のためのデータ分析



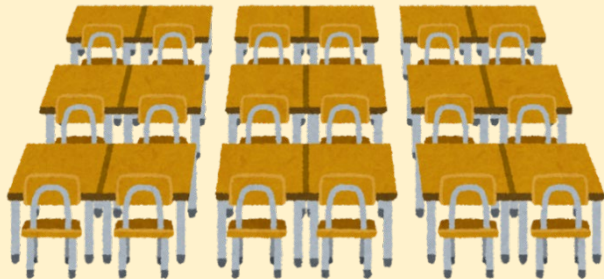
## 公開・共有

- 他の研究者の二次分析に資するデータとプログラムの公開



## 教育・学習

- 学生たちにデータ分析をさせるゼミ・講義・演習など

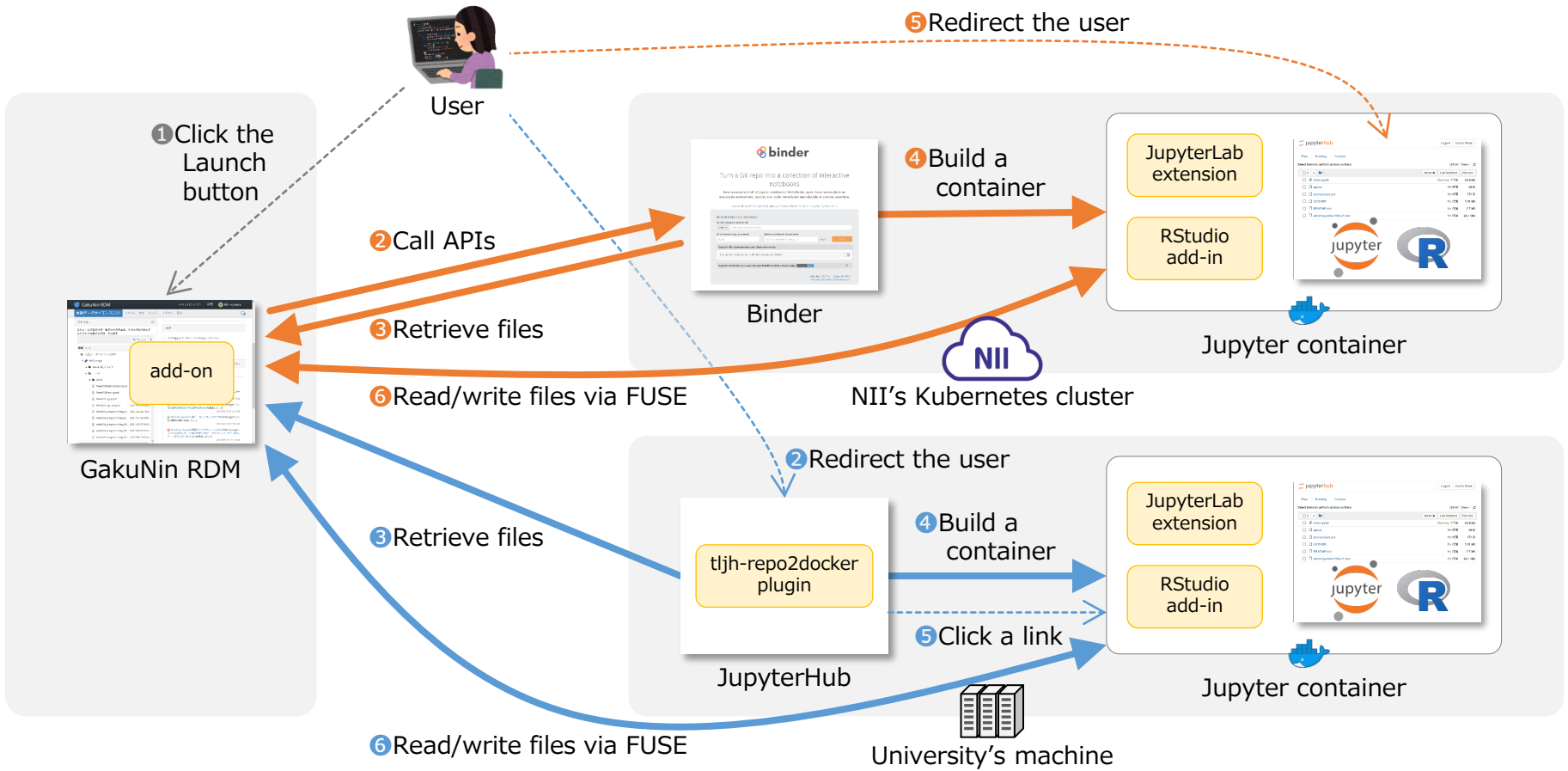


## 引き継ぎ

- 先輩の研究環境を後輩が再現し、研究を継続する



# システム構成



# 類似サービスとの比較

	GakuNin RDM データ解析機能	GESIS Notebooks	mybinder.org	Google Colab	Microsoft Codalab
対象分野	汎用	社会科学	汎用	主に深層学習	深層学習
提供元	NII (日・学術機関)	GESIS (独・学術機関)	Project Jupyter (任意団体)	Google (米・民間企業)	Microsoft (米・民間企業)
アカウント	学認	GESIS	不要	Google	Codalab
対応言語	Python, R, Shiny, MATLAB (準備中)	Python, R, Julia	Python, R, Julia	Python, R, Julia, Swift	Python
対応リポジトリ	GakuNin RDM, JDCat (WEKO3), GitHub, Gist, GitLab, Zenodo, Figshare, Hydroshare, Dataverse	GitHub, Gist, GitLab, Zenodo, Figshare, Hydroshare, Dataverse	GitHub, Gist, GitLab, Zenodo, Figshare, Hydroshare, Dataverse	Google Drive, GitHub	?
メモリ CPU ストレージ	3GB 36コア共有 10GB	32GB 2コア 10GB	2GB 1コア ?	13GB 2コア 40GB	?
タイムアウト	半永続 30日不使用で消去	40分	10分	90分 / 12時間	?
インフラ	オンプレ	オンプレ	Google, OVH, Turing Institute	Google	Microsoft

# 詳しい情報

---

## 誰が使えますか？

- 九州大学のユーザーは全員、GakuNin RDM データ解析機能をご利用いただけます。

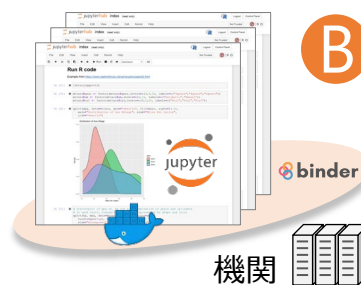
## どこを見れば分かりますか？

- マニュアル <https://support.rdm.nii.ac.jp/>  
→ ユーザーマニュアル → データ解析機能
- 解説動画 [https://youtu.be/\\_FzOpDTQrBQ](https://youtu.be/_FzOpDTQrBQ)

## 分からないときは誰に聞けばいいですか？

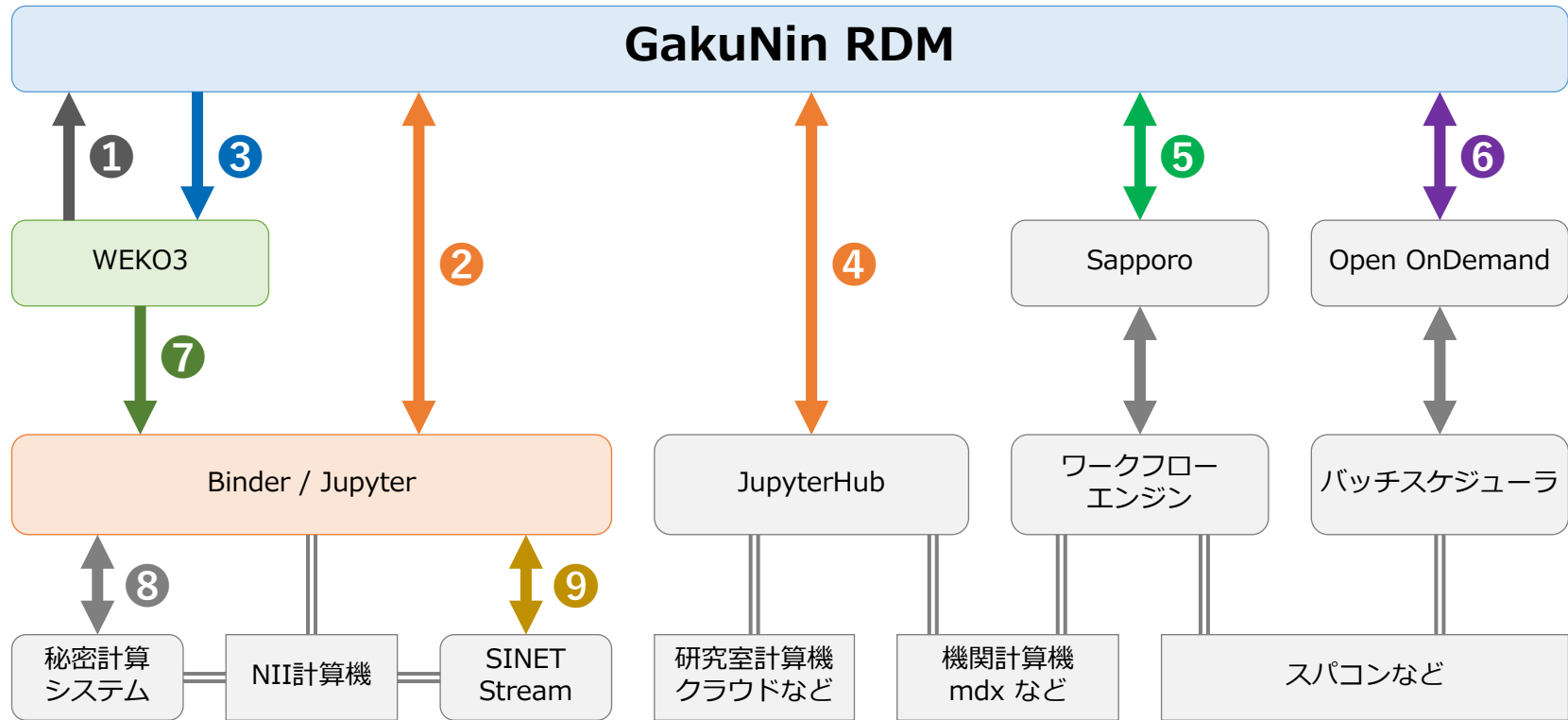
- 九州大学 DX推進本部 研究データ管理支援部門  
[rds\\_help@dx.kyushu-u.ac.jp](mailto:rds_help@dx.kyushu-u.ac.jp)

# 外部計算機連携のバリエーション



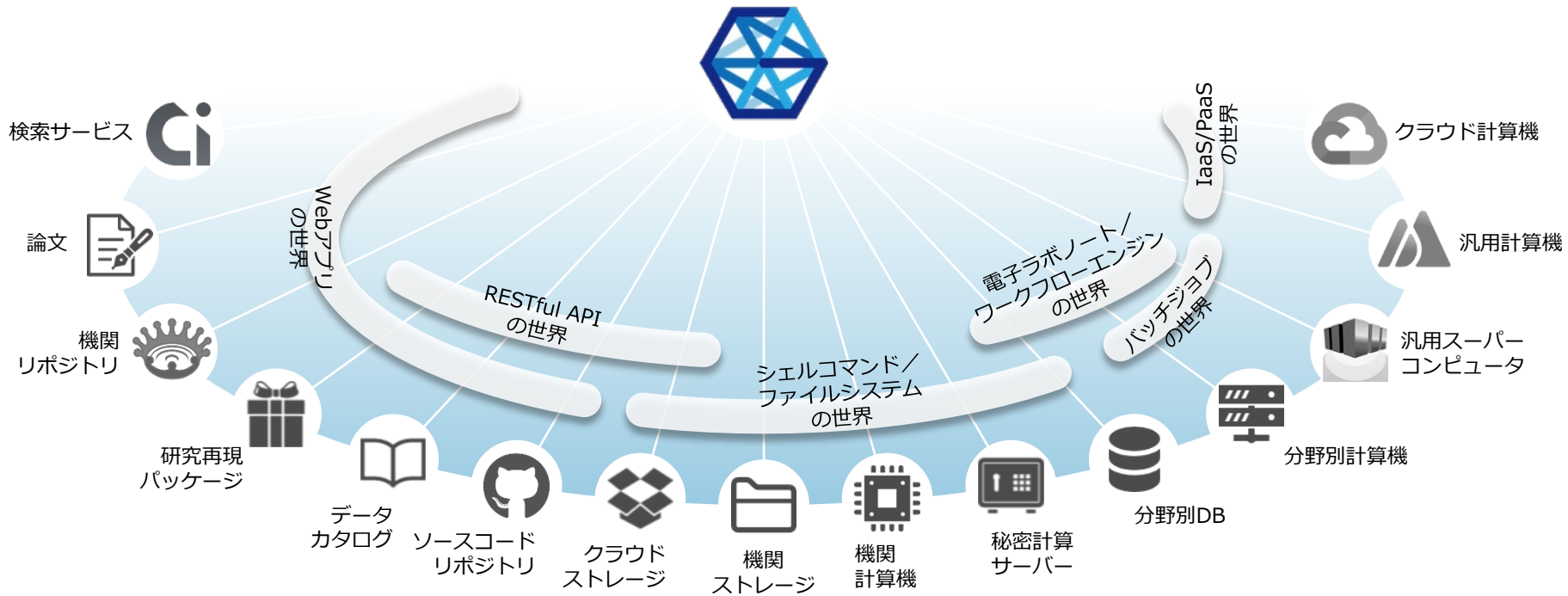
	システム	運用主体	認証方法	同時起動数	ドメイン名+サーバ証明書	バックエンド
A	Binder	NII	学認	10個/ユーザ	○	Kubernetes
B	Binder	機関	学認, OAuth, LDAP, etc.	任意設定	必要	Kubernetes
C	JupyterHub	研究室等	OAuth, LDAP, ローカル	1個/ユーザ	不要	Linux VM

# コード付帯機能群



①③ 計算再現パッケージ機能	GRDMプロジェクトをWEKOで公開、他者がGRDMに取り込み再利用	開発中
②④ GakuNin RDMデータ解析機能	Jupyterによるデータ解析環境をGRDMから構築	運用中
⑤ 外部ワークフローエンジン連携機能	Sapporo 経由でワークフローを実行、結果をGRDMに回収	設計中
⑥ 外部HPC連携機能	Open OnDemand からGRDMのファイルを読み書き	提供中
⑦ WEKOオンライン分析機能	Jupyterによるデータ解析環境をWEKOから構築	運用中
⑧ 秘密計算システム統合機能	秘密分散によるセキュアな解析環境をJupyterから利用	設計中
⑨ SINETStream連携検討	SINETStreamによるリアルタイムデータ収集環境を構築	開発中

# 目指したい将来像



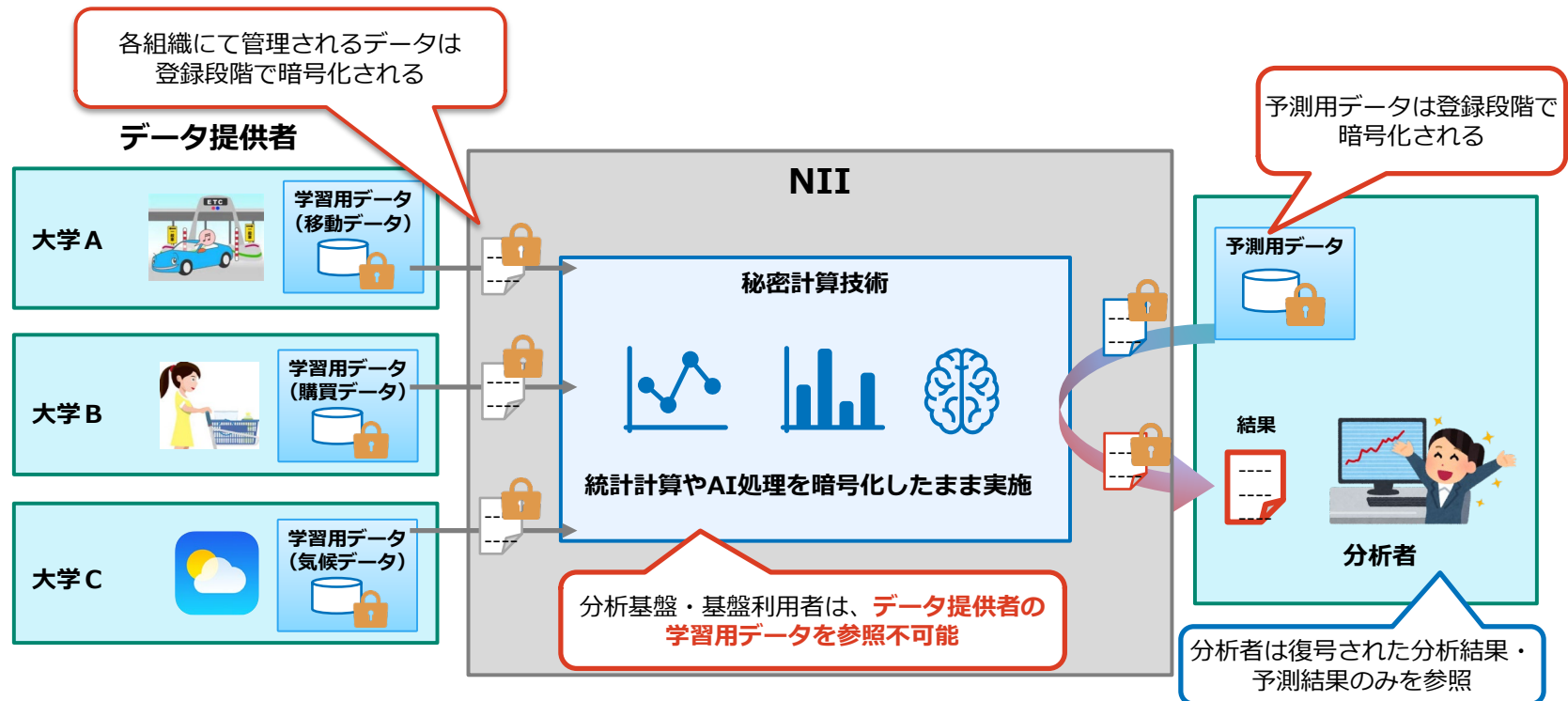
**GakuNin RDM を核として  
データの世界と計算機の世界を結ぶ**

**RCOS**  
rcos@nii.ac.jp



# 秘密計算システムの大学向けトライアル

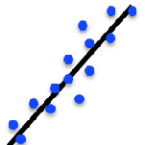
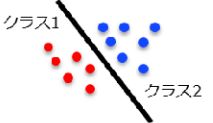
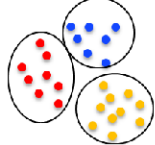
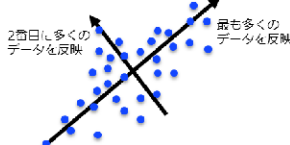
- NIIの計算機上にNTTの秘密計算システムを構築し、大学等の研究者に無償提供。現在、実験パートナーを募集中
- 研究者は、システムに手持ちのデータを登録し、秘密計算ライブラリを用いたデータ分析プログラムを書いて実行
- 秘密分散に基づくマルチパーティ計算の実用性を評価していただく → 有用な秘匿解析機能の開発へ



# 秘密計算システム「算師」

## 秘密計算AIで実現したAI 4大カテゴリの 主要なアルゴリズム



①回帰（連続値の推定）	②クラス分類（種類毎の分類）	③クラスタリング（類似値による分類）	④データ次元圧縮（主成分の抽出）
教師あり学習		教師なし学習	
			
<p><b>ニューラルネットワーク (FFNN、CNN、RNN)</b></p> <ul style="list-style-type: none"> <li>データ量：大</li> <li>特徴パラメータ量：極めて大</li> <li>予想外の発見的な結果を得やすい（ビッグデータの共有による解析精度向上）</li> </ul> <div style="display: flex; justify-content: space-around;"> <div data-bbox="154 821 473 892">街の店舗横断の購買予測</div> <div data-bbox="521 821 840 892">複数の病院横断の血液分析による疾患予測</div> </div>		<p><b>階層型クラスタリング k-means</b></p> <ul style="list-style-type: none"> <li>データ間の類似度にもとづいてデータをグループ分けする手法</li> </ul> <div style="display: flex; justify-content: space-around;"> <div data-bbox="1014 878 1362 949">複数の病院横断の疾病患者グループ分析</div> <div data-bbox="1014 992 1362 1063">街の店舗横断の顧客グループ分析</div> </div>	
<p><b>決定木、GBDT</b></p> <ul style="list-style-type: none"> <li>データ量が小～中</li> <li>特徴パラメータ量 小～大</li> <li>データが少量でも精度を出しやすい（データ数が揃わない時の分析精度向上）</li> </ul> <div style="display: flex; justify-content: space-around;"> <div data-bbox="154 1056 473 1128">複数の病院横断の希少疾患予測</div> <div data-bbox="521 1056 840 1128">クレジットカード会社横断の不正取引検知</div> </div>		<p><b>主成分分析</b></p> <ul style="list-style-type: none"> <li>多量の変数を少量の変数に置換、要約</li> <li>総合力に影響している項目が把握できることで分析を効率化する</li> </ul> <div style="display: flex; justify-content: space-around;"> <div data-bbox="1497 878 1816 949">街の店舗横断の購買傾向分析</div> <div data-bbox="1497 992 1816 1063">製造メーカー横断の製品評価分析</div> </div>	
<p><b>Lasso回帰</b></p>	<p><b>ロジスティック回帰</b></p>		

# 詳しい情報

---

## お申し込み方法

- 以下の内容を記載したメールをお送りください。
  - 研究者の氏名、所属、連絡先メールアドレス
  - 研究内容の概略、秘密計算を利用するメリット
- 個別にヒアリングを行い、詳細を決定します。
  - 原則として1チームにつき1セットの秘密計算環境を提供します。

## 資料

- トライアル概要  
<https://rcos.nii.ac.jp/news/2023/01/20230123-0/>
- ニュースリリース  
<https://www.nii.ac.jp/news/release/2023/0123.html>

## お問い合わせ

- [sc-trial2023@nii.ac.jp](mailto:sc-trial2023@nii.ac.jp)