

第3回データワークショップ

SPring-8データセンター構想の進捗報告

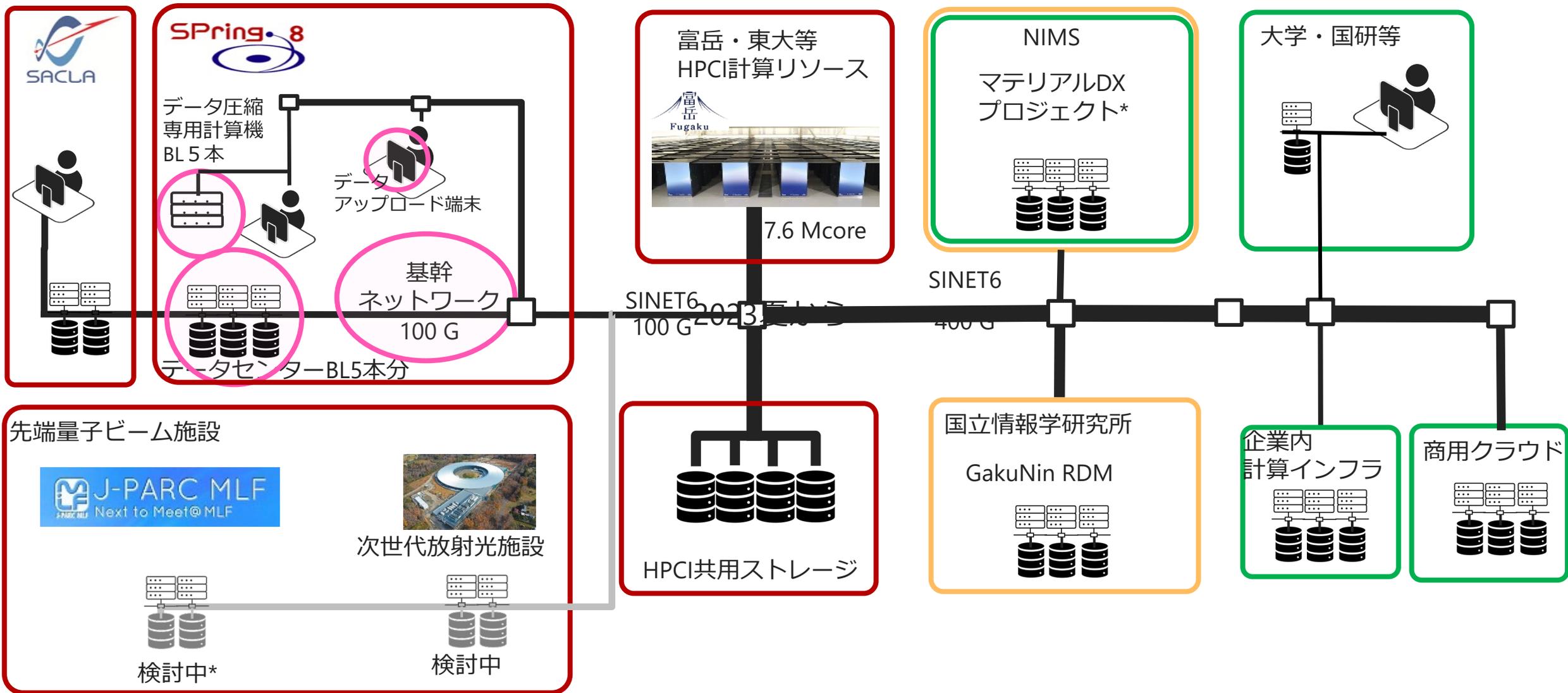
初井宇記

理化学研究所

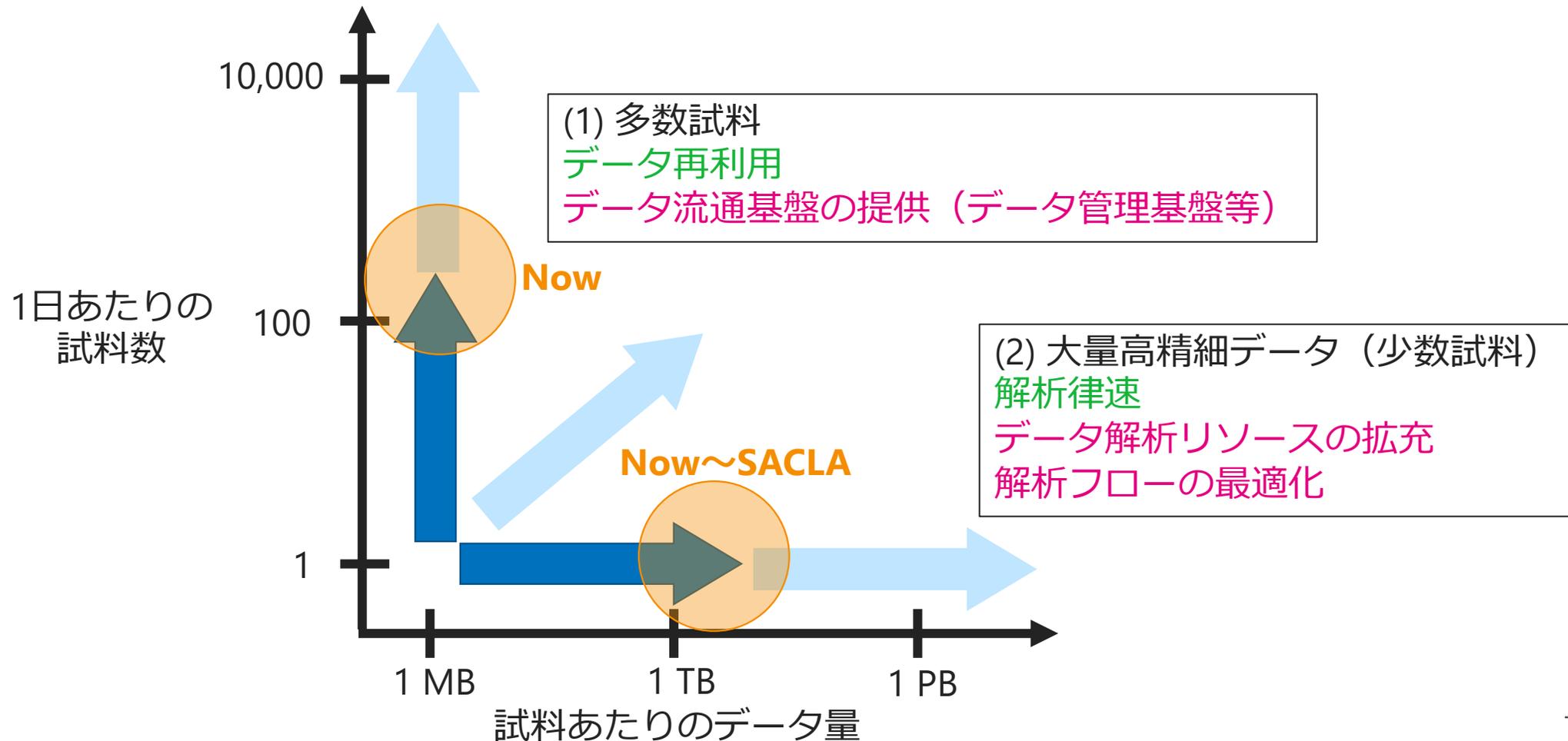
アウトライン

- SPring-8データセンター構想の進捗報告
 - 類型化と課題、検討中の機能
- スケジュール
- まとめ

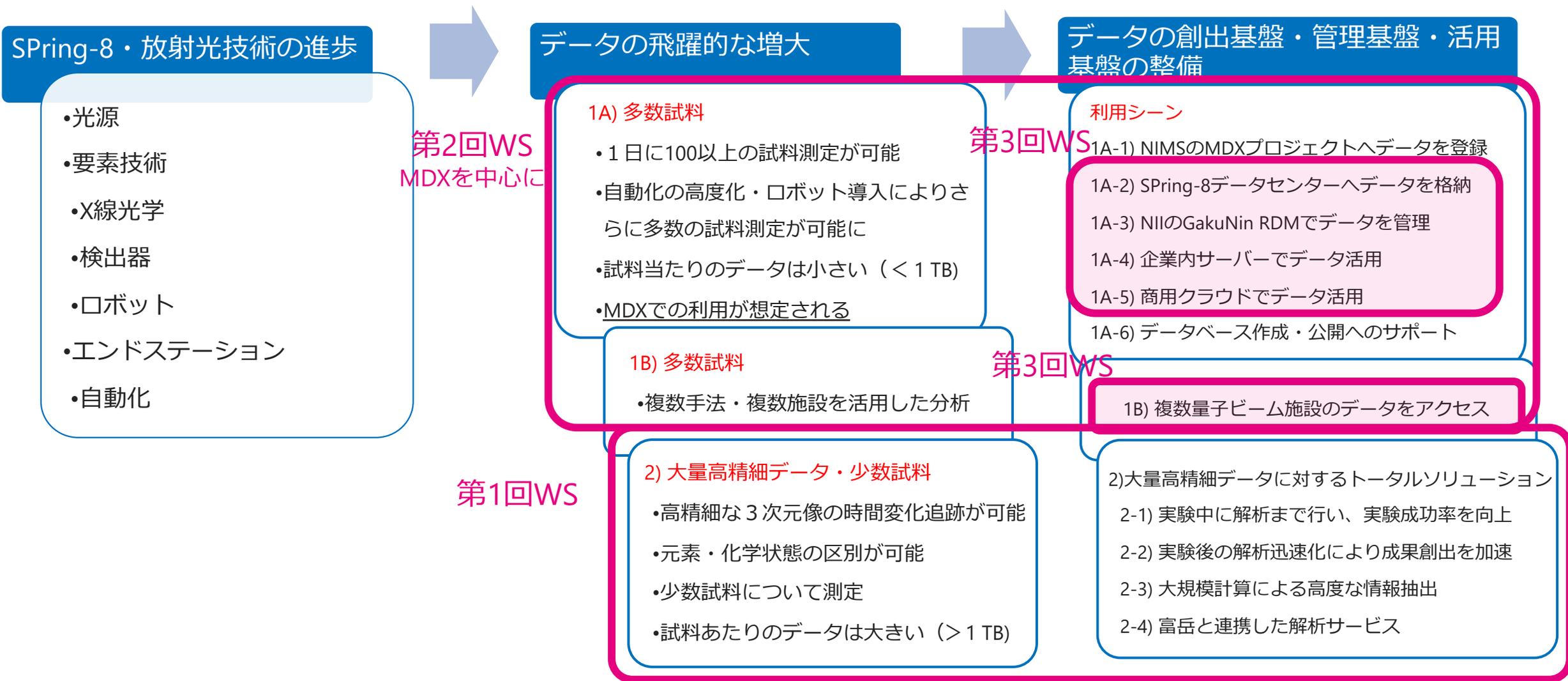
SPring-8データセンター構想: インフラ



* クラウド上に設置を予定



SPring-8データセンター構想の動機・現状および検討中のサービス



1A-2) SPring-8データセンターへデータを登録

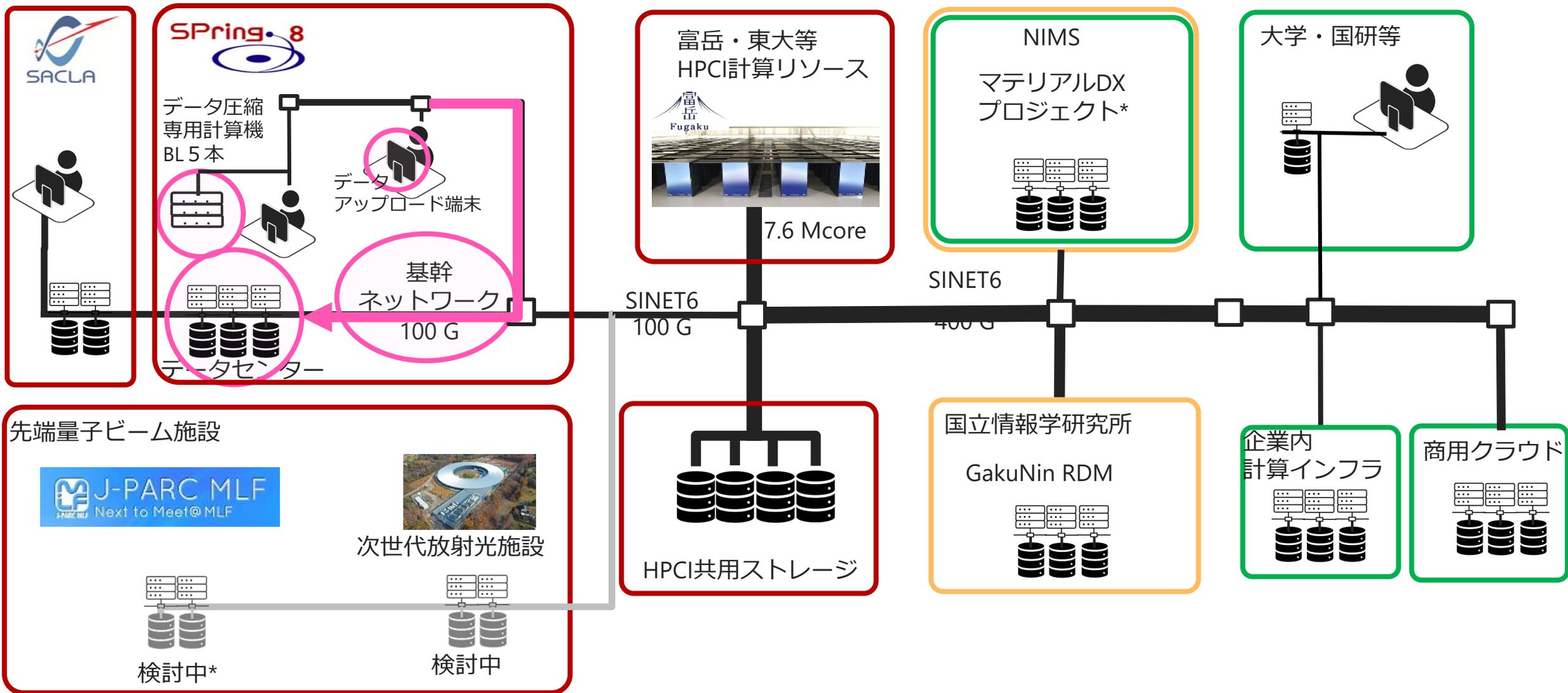


ビームラインに設置している端末

- SPring-8データセンターのアップロード用のweb画面が表示される。
- 自動化されたビームラインは半自動でアップロードされる。

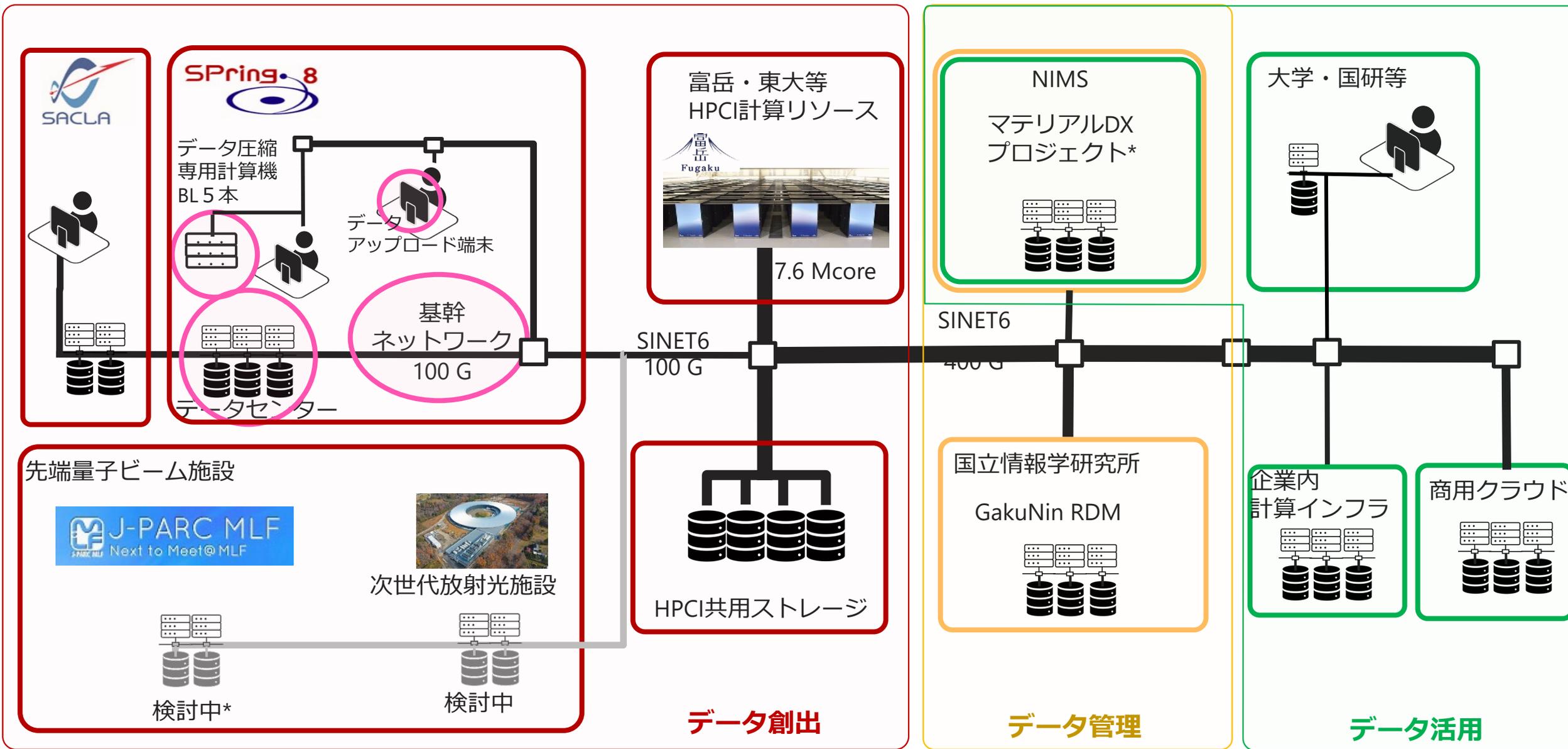
自動化ビームラインについて具体的に検討している例については次の平木の講演をお聞きください。

SPring-8データセンター構想: インフラ



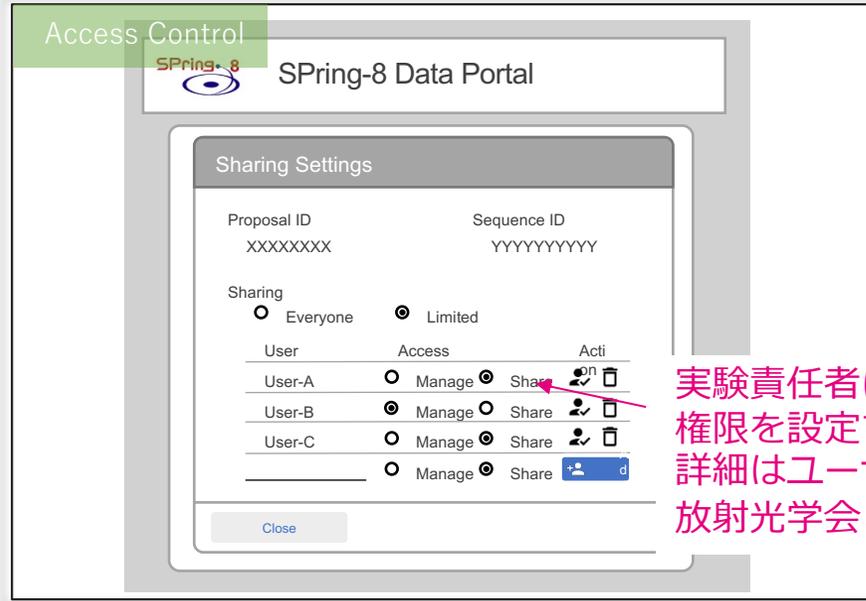
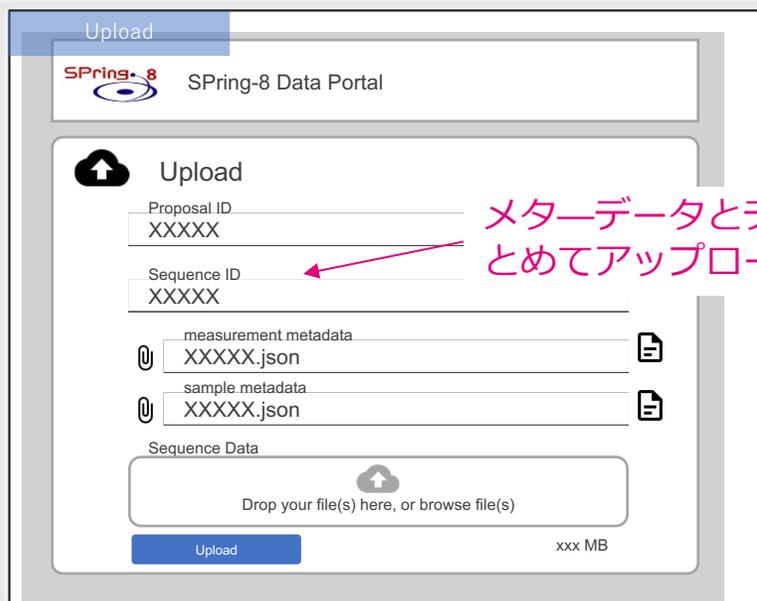
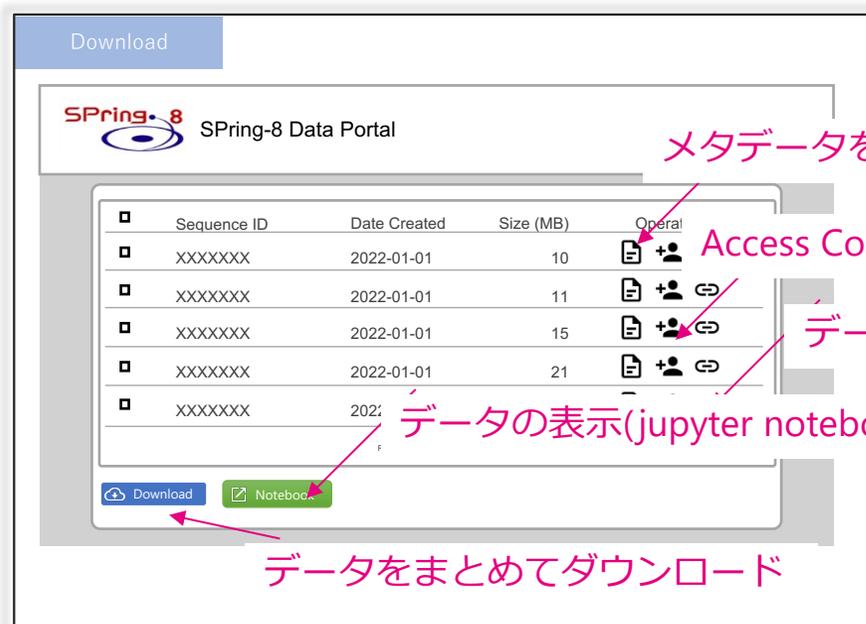
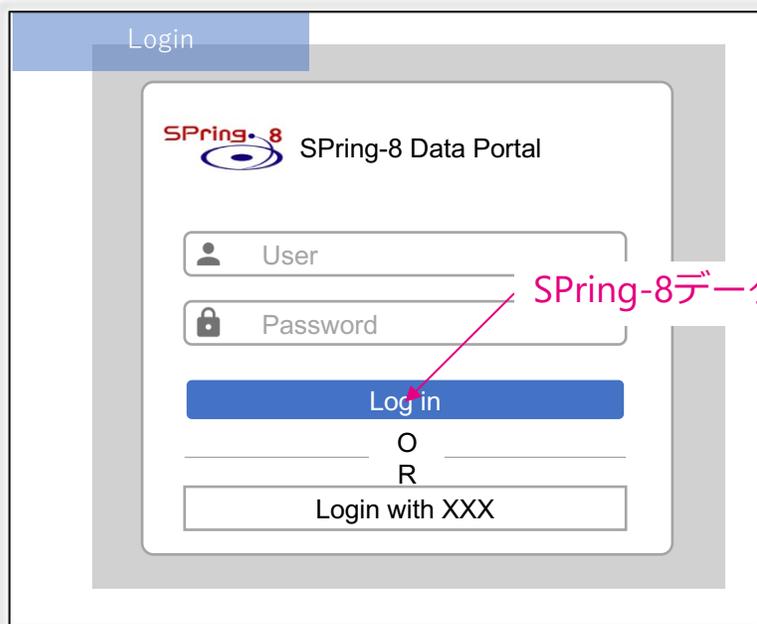
* クラウド上に設置を予定

SPring-8データセンター構想: インフラ



* クラウド上に設置を予定

SPring-8 Data-Flow Systemのユーザー画面 (検討中の案)



データの管理について

課題番号 XXXXXXXXX

Sequence ID XXXXXX

Sequence ID XXXXXX

Sequence ID XXXXXX



研究者A



研究者B

- 課題番号

- (実装案) 課題番号毎にアクセス権限を設定できるようにする
 - 皆様のご意見をお待ちしております。
 - Sequence毎のアクセス管理が必要な場合はお知らせください。

- Sequence

- 一定の時間測定したデータをsequenceとよび、メタデータを付与する。
- 以下のメタデータ付与は義務とする。
 - 管理用メタデータ (課題番号、実験日など)

メタデータ

第2回データワークショップで提示させていただいた素案

前回からの進捗

SPring-8 Data-Flow System メタデータは以下の3種とする

- 管理情報
- 測定関係情報
- 試料関係情報

背景

• 理研では統合情報本部においてデータ構造化を進めている
現状

- SPring-8についても連携してメタデータ構造を改訂中
- 世界標準の規格をなるべく参照
- NIMSのデータ構造と連携できる構造化とする
- 放射光学会でのデータ構造化の議論とも整合性のある形にする

SPring-8 Data-Flow Systemでのメタデータの入力形式 (素案)

稼働から五年後にユーザー数1万5000人、
データ数1億程度と想定し検討中

Spring-8データセンター構想における
metadataの素案

2022/2/24
本村幸治
理化学研究所
放射光センター

改訂中

Metadata

- 目的
- 本稿の前提
- 方法

- 理研 本村中心に分野非依存の共通部分から技術検討を進めている。
ご協力・ご支援いただける方は是非ご連絡をお願いします。
現時点の概要はワークショップのweb pageにアップロードしています。



理化学研究所
情報統合本部
小林 紀郎 (D.Eng.)

是非インプットをお願いします。

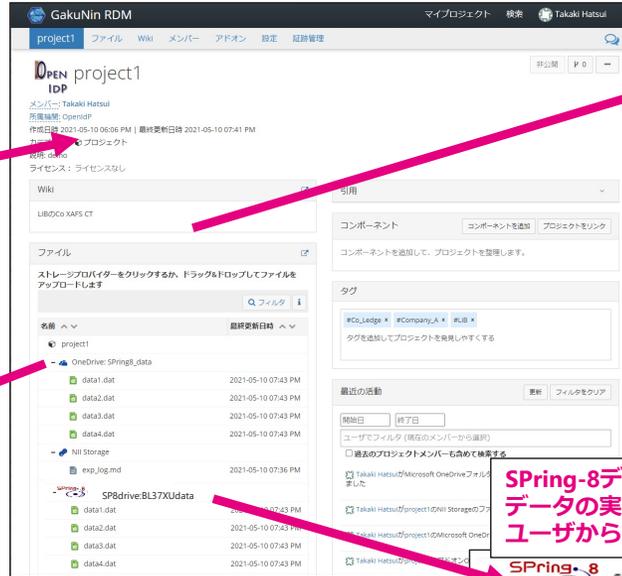
1A-3) NIIのGakuNin RDMでデータを管理

NII上に構築されているGakunin RDM or プロジェクトのデータ管理基盤



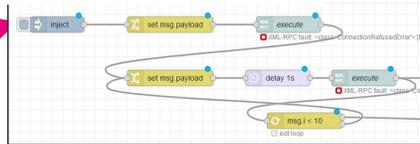
所属機関サーバ
プロジェクトサーバ
商用クラウド

研究テーマごとのページ



SPring-8データセンター or HPCI共用ストレージ上 (新規) 実験中の解析を最先端のアルゴリズムで実施できる (ユーザのブラウザ上)

フローベース: 一般ユーザ向け



コマンドベース: 専門家向け

```
*** This software and the associated documentation are confidential ***  
*** and proprietary to Synopsys, Inc. Your use or disclosure of this ***  
*** software is subject to the terms and conditions of a written ***  
*** license agreement between you, or your company, and Synopsys, Inc. ***  
[ui000002@login1 ~]$
```

SPring-8データセンター or HPCI共用ストレージ上 (新規) データの実体はSPring-8が管理するサーバにあるが ユーザからは手元に見える

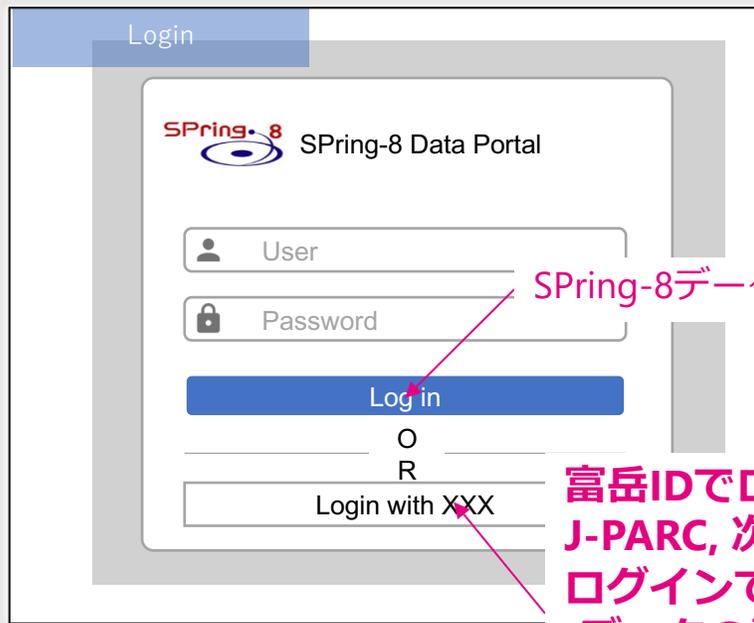


研究室の他の装置のデータと一括管理が可能
SPring-8のデータは実体はSPring-8にあるが参照先が分かる。
解析もここからアクセス可能とする予定

データ公開機能もあり

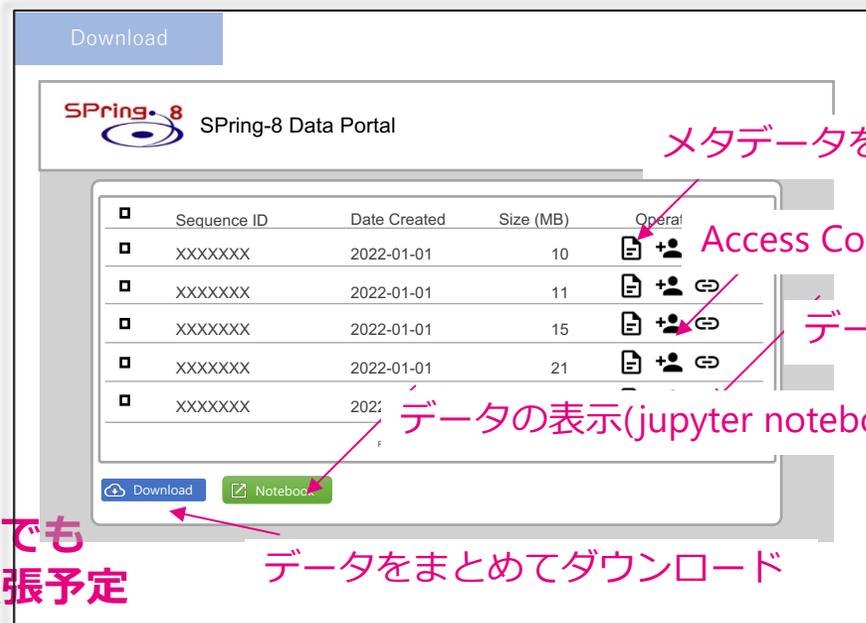
国立情報学研究所(NII)にはWEKO3というサービスがあり、この機能を使えば GakuninRDMで管理されているデータを公開することも可能。

1B-1) 複数量子ビーム施設のデータをアクセス



SPring-8データセンターのID

富岳IDでログイン
J-PARC, 次世代放射光IDでも
ログインできるように拡張予定
(データの認証・認可
先端大型施設間で連携)



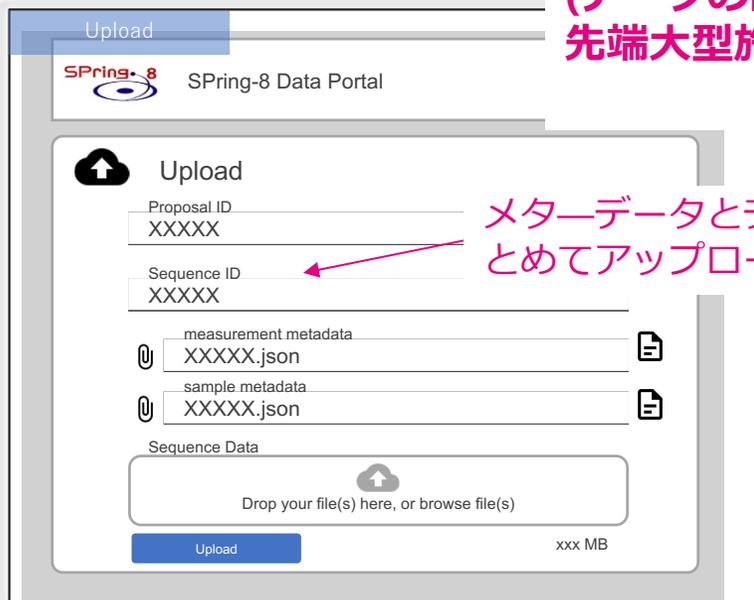
メタデータを閲覧

Access Controlへのリンク

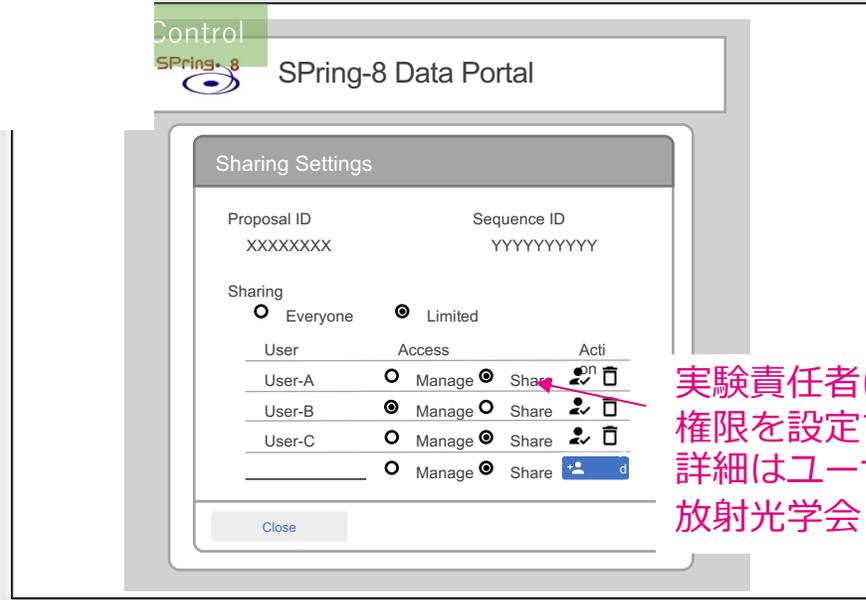
データへのリンク(REST API)

データの表示(jupyter notebook)

データをまとめてダウンロード



メタデータとデータをまとめてアップロード

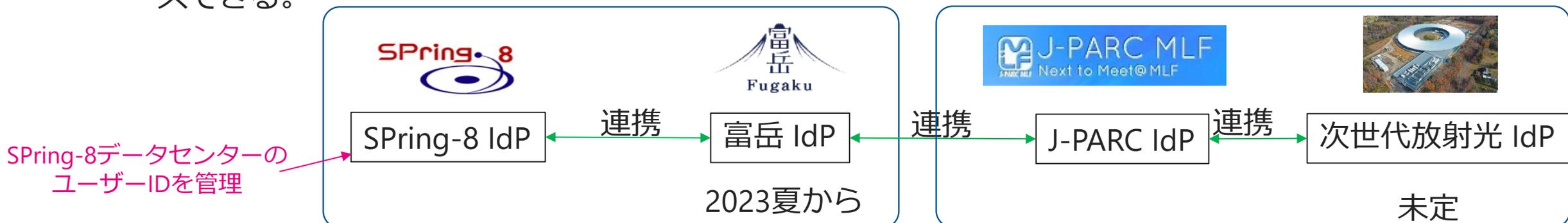


実験責任者はデータのアクセス
権限を設定することができる。
詳細はユーザーコミュニティ・
放射光学会と協議

1B-1) 複数量子ビーム施設のデータをアクセス

- 連携が必要な機能

- データの認証・認可：共通化することで一つのプログラムから複数施設のデータを簡単にアクセスできる。



将来的にNIIの認証基盤とも連携

- 共通化が必要な機能

- データアクセス方法の共通化
例：データアクセスの基幹部分のソフトウェアがHOST名を変えるだけでどの施設でも使える

検討中の実装

Amazon社のオブジェクトデータサービスS3と互換性のあるサービス

ユーザーから見た使い方

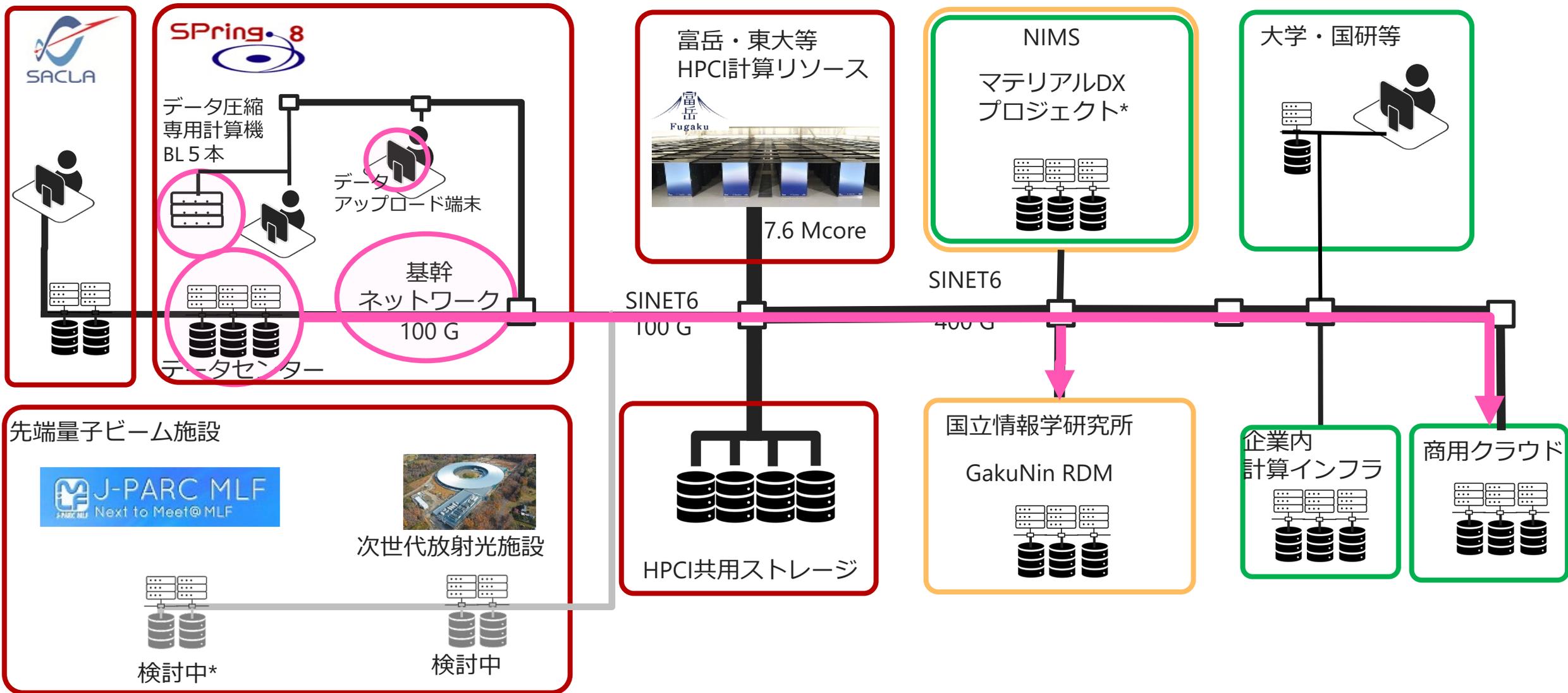
pythonライブラリ boto3によりSPring-8データセンター内のデータをアクセスできるようになる。

1A-4) 企業内サーバーでデータ活用

1A-5) 商用クラウドでデータ活用

- ヒアリング状況
 - 企業内サーバへのデータ移動
 - VPNなどを介するため帯域が狭く、実務上ダウンロードできない場合が多い。
 - 商用クラウド (Amazon AWS, Microsoft Azure他)
 - 利用できる企業が増えている。
- サービス方針案
 - 1A-4) 企業内サーバーでデータ活用
 - 小さなデータ (<< 1 GB)の場合ダウンロードできる
 - SPring-8 Data-Flow Systemで対応
 - 1A-5) 商用クラウドでデータ活用
 - SINET6に直接接続しているデータセンターを対象に、中程度のデータでもダウンロード可能なサービスを検討する。
 - 対応できるデータサイズはこれから性能確認を行う
 - SPring-8 Data-Flow Systemで対応
 - プログラムから直接読み出すユーザにはboto3によるデータアクセスを提供

SPring-8データセンター構想: インフラ

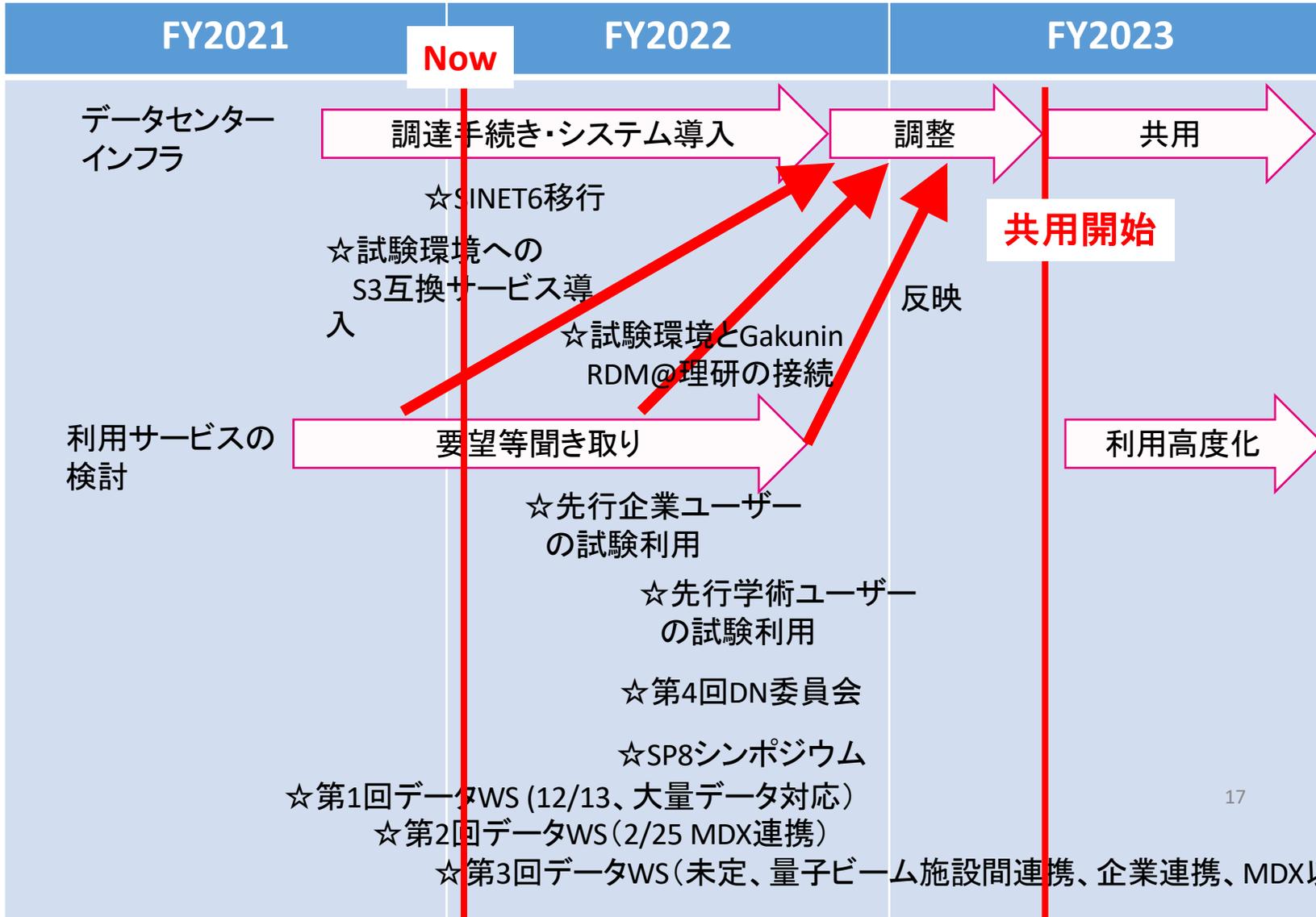


* クラウド上に設置を予定

新設

2021年度補正予算

スケジュール



- MDXプロジェクト参画代表者
 - 専用ビームライン2機関
 - 企業 2法人
- などからフィードバック・ご要望をいただいています。

17

まとめ

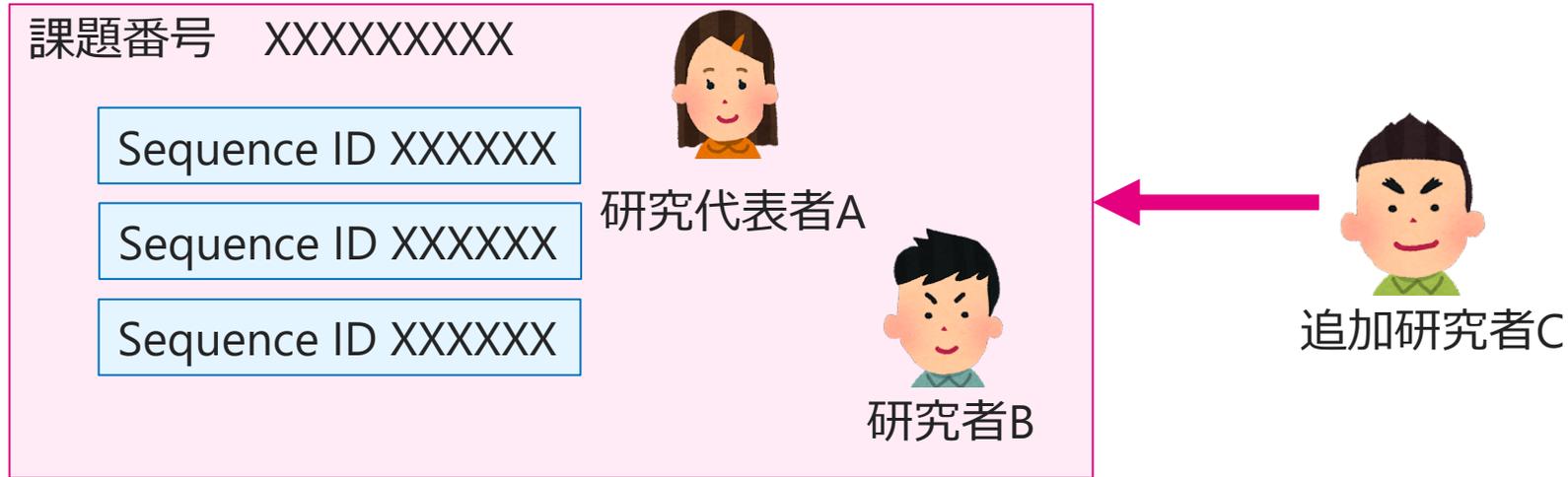
- SPring-8データセンター構想の進捗報告
 - 類型化と課題、検討中の機能
 - データ管理の考え方
 - 課題番号：アクセス管理するレイヤー
 - Sequence: SPring-8データセンターでデータを管理するに当たり、新たに導入するデータの粒度
- スケジュール
- SPring-8データセンターに関する情報
 - <https://dncom.spring8.or.jp/>
- ご意見のある方は、是非ご連絡をお願いします。
 - dncs@spring8.or.jp

謝辞

- [SPring-8データセンター構想]
 - 城地保昌(JASRI,理研),中嶋享(JASRI), 本村幸治(理研), 平木俊幸(理研),杉本崇(JASRI),
 - 山鹿光裕(JASRI), 矢橋牧名(理研)
 - 中町将貴 (理研) 、 渡邊一輝 (理研)
- [HPCI連携、データ圧縮、Workflowツール、認証]
 - 理研計算科学研究センター(R-CCS, 富岳)
 - 松田元彦, 原田浩、金山秀智, 山本啓二, 佐藤賢斗,
 - 佐野健太郎, 庄司文由, 松岡聡
- [研究データ管理]
 - 理研情報統合本部(R-IH)
 - 實本英之, 富士健太郎,小林紀郎, 美濃導彦
- 国立情報学研究所(NII)
 - 込山悠介, 山地一禎, 喜連川優

SPARE SLIDE

データの管理について



- 課題番号

- (実装案) 課題番号毎にアクセス権限を設定できるようにする
- 皆様のご意見をお待ちしております。
- Sequence毎のアクセス管理が必要な場合はお知らせください。

- Sequence

- 一定の時間測定したデータをsequenceとよび、メタデータを付与する。
- 以下のメタデータ付与は義務とする。
- 管理用メタデータ (課題番号、実験日など)

データのオープン化と SPring-8データセンターの関係

1A-6) データベース作成・公開へのサポート

Big Data Value Chain[1]



[1] Freitas A., Curry E. (2016) Big Data Curation. In: Cavanillas J., Curry E., Wahlster W. (eds) New Horizons for a Data-Driven Economy. Springer, Cham. https://doi.org/10.1007/978-3-319-21569-3_6

[2] digital object identifier

SPring-8の場合



SPring-8データセンター構想

- 解析律速の解消
SPring-8内データセンターおよび富岳等HPCIスパコンによる計算リソース
- データ共有の簡便化
SPring-8 data-flow system

ユーザーコミュニティやプロジェクト（ARIM、MDX等）が主導

データ公開のために必要なインフラ：永続的なデータ保存、検索・アクセス基盤の整備

データ公開基盤（各研究機関、国立情報学研究所、NIMS等）を担当される機関と連携し、データ公開化に協力して参ります。

データのオープン化に関するSPring-8データセンターの役割

考え方

SPring-8データセンター構想

連携機関と協力してSPring-8ユーザーにデータに関する課題を解決していく

SPring-8に設置するデータインフラは以下を担う

データ創出

実験中のデータ処理

一時保管(論文出版に必要な期間を目標としています。)

データの所有権、使用权

データネットワーク委員会にて詳細を決定して参ります。

ご意見がある方は、データネットワーク委員会の委員、もしくは理研側へご連絡ください。

前提としては、以下で進めるものと考えています。

【所有権】データの所有権はユーザが持つ

【使用权】施設者は保守・高性能化などに資する情報を得るためにデータを使用する場合がある。

ただし、収集したデータは第3者へは開示しない

【運用】データはバックアップをとらない。ユーザーは期限が来たら消去する必要がある。