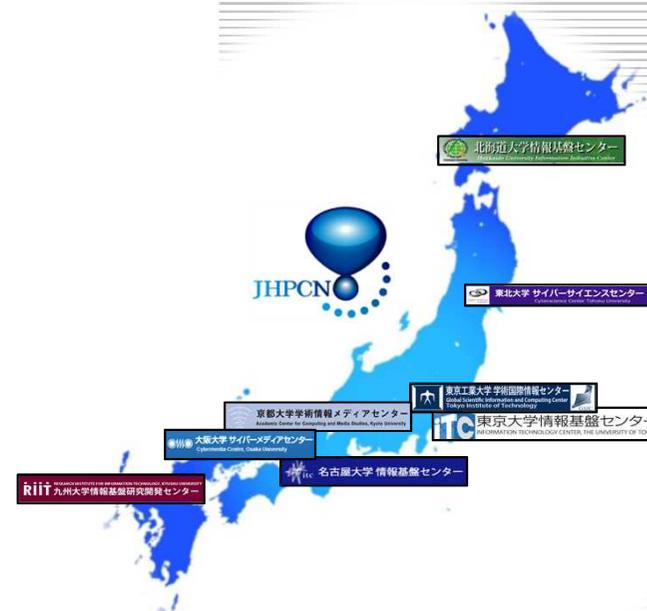


情報基盤センター群の取り組みと データ活用社会創成プラットフォームについて

田浦 健次郎
東京大学 情報基盤センター



自己紹介

- 東京大学
- 本務大学院 情報理工学系研究科
- 兼務 情報基盤センター長（2018.4～）
 - JHPCN（共同利用共同研究拠点）統括拠点長
 - HPCI（コンソーシアム理事、計画推進委員会、第2階層まとめ役、etc.）
 - mdx仕様策定委員長
- 研究分野
 - システムソフトウェア（とくに並列処理、並列プログラミング処理系、データ処理系）

次期学術情報基盤に求められること

第6期科学技術・イノベーション基本計画

- 新たな研究システムの構築（オープンサイエンスとデータ駆動型研究などの推進）
 1. 信頼性のある研究データの適切な管理・利活用促進のための環境整備
 2. 研究DXを支えるインフラ整備と高付加価値な研究の加速
 3. 研究DXが開拓する新しい研究コミュニティ・環境の醸成

学術・産業・社会ニーズ

データ駆動科学
大規模データ処理・ML・AI

Society 5.0, DX

オープンサイエンス、データ公開

研究データ管理・不正防止

教育の進化・データ利用・LA・教育デジタルコンテンツ利用

直面している課題・問題

多くの分野で必要な計算・データ基盤が拡大し、分野ごとの整備は困難・非効率に

分野ごとのデータ基盤整備計画はあるが、しばしば計算（解析・利活用）能力が足りない

データ科学の需要は従来のHPC用環境（スパコン）だけでは満たせない

パブリッククラウド環境は高価かつほぼ海外

オンプレクラウドの導入は技術的なハードルが高い（&国内ベンダーの蓄積が少ない）

- **高性能データ&計算基盤**
- **分野を超えた共通基盤**
- 新しい学際的連携を生む**共同利用共同研究プログラム**（コミュニティ創成）
- 研究データ基盤（**NII Research Data Cloud**）含め、システム間「連携」を前提とした設計

データ科学のインフラ具体像

- 機械学習フレームワーク
- JupyterなどWebベースインタラクティブ環境
- 分野ごとに異なるアプリケーションソフトウェア (流体、材料、ゲノム、etc.)
- 分野用データレポジトリなどの常駐サービス

- フィールドのIoTデバイス (センサなど) からのデータ収集
- キャンパスの実験機器などからのデータ収集
- キャンパス・サイトをまたがったデータ処理
- 研究室データの処理
- リモートデスクトップなどGUI環境

要求の発生場面

- 機密性の高いデータの蓄積・処理
- 個人に由来するデータの蓄積・処理

- 探索的データ処理
- 複数・多数のデータの活用

システム要件

1. ユーザごと、分野ごとに柔軟に構成可能な環境
2. 自由度の高い構成が可能なネットワーク環境
3. ユーザごと、分野ごとに分離・隔離されたセキュアな環境
4. データ検索・発見・共有から計算へシームレスに移行できる環境

• ≈ 仮想化されたマルチテナント (クラウド) 環境

- データを蓄積・共有できるだけでは足りない
 - 高性能計算機と一体化の必要 (特にML・AI利用)
- これまでの高性能計算機だけでも足りない
 - GPUとMLフレームワークがあれば済むわけではない

データ活用社会創成プラットフォームとは

● 第2回情報委員会(令和1年8/7) 資料

https://www.mext.go.jp/content/20200116-mext_jyohoka01-000004142_3.pdf

**Society5.0を実現するためのデータ活用による知識集約型社会の創成
ーデータ活用社会創成プラットフォームの構築ー**

データ活用社会における現状認識

- ICT機器の爆発的な普及や、AI、ビッグデータ、IoT等の社会実装が進むなど競争が激化、一部の企業や国のデータの囲い込みにより経済社会システムの健全な発展が阻害される懸念。
- 我が国が成長していくためには、デジタル新時代において、データを我が国全体の共同資産として、スピード感をもったデータ活用環境の整備が急務。
- Society5.0が目指すインクルーシブな社会を実現するためには、地域における知識集約の中核を担う大学を起点としてイノベーションの創出を図り、知識集約型社会を構築することが重要。
- サイバー空間とフィジカル空間が融合するデジタル新時代において、我が国に蓄積された農業、医療・健康の分野、教育データを含む多くの有用なビッグデータを共同で活用する上で、人材と技術を有する全国の大学を超高速・高信頼で網目状につなぐ国際的優位性をもつSINETを最大限活用することが重要。
- 異種データや異種知識の融合・活用を促進するための「場」として、様々な分野のデータ保持者、解析者、利用者が参画するコミュニティを形成するとともに、データ活用を目指す利用者へのコンサルティングやアプリケーション開発支援が不可欠。

文部科学省における取組

- 経済財政運営と改革の基本方針2018や未来投資戦略2018において重要性が指摘されているリアルデータの利活用を念頭に、データ活用社会創成プラットフォームを推進するため、SINETを通じて収集されるリアルデータの集積や、解析結果を速やかにフィードバックする機能を備えたシステムを整備(2019年度予算)。
- 文部科学省と大学コミュニティ、地域社会等が一体的に連携し、全国の国立大学等をハブとしたデータ活用社会創成プラットフォームの実現促進に向けた検討を行うための「データ活用社会創成プラットフォームの推進に関する有識者会合」を設置。
- 地域・産業・社会基盤を支える拠点となる大学を中心として、民間への利用拡大も視野に我が国全体の知識集約型社会の実現に向けた環境「データ活用社会創成プラットフォーム」を構築。

データの高度利用環境(NII・東大に先行して整備)
【設備整備】
IoT接続(モバイル)
AI特化PC
リアルタイム処理対応サーバー
高速/セキュリティ等

SINETを通じて、全国のデータ収集・通信・解析環境をオンデマンドで活用。
高度・多様なデータ利活用により新たな価値を創出。

データ活用社会創成プラットフォームの推進に関する有識者会合
リアルタイム処理対応基盤社会創成プラットフォームの実現に向けた実務的な検討を行う場
【主な検討課題】
・リアルタイムデータの解析・活用を目的とした基盤ソフトウェアの研究開発や技術の実証のための基盤システムの整備のあり方
・産学連携体制(コミュニティ)の構築・強化、その中核としての大学の役割等一体的な連携を確保する仕組み 等

文部科学省と大学コミュニティ、地域社会等が一体的に連携し、プラットフォームの実現に向けて整備・検討を加速

大学等連携コンソーシアム
大学を中核としたデータ活用実務機関が連合したコンソーシアム
【主な取組】
・データプラットフォームの活用促進、データ活用ニーズ調査
・コミュニティ間連携の強化・促進等

活用ニーズを踏まえたシステム整備・ソフトウェア開発
【大学等におけるデータ利活用の潜在的なニーズ】
・地域農業・漁業・観光業のスマート化
・認知症・生活習慣病などの早期発見、予防方法の提案
・スポーツ科学への応用
・初中段階から高等教育、社会人教育に至る一貫した教育データの利用 等

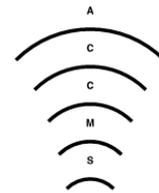
有識者会合（主査：安浦先生）

基盤の先行整備

...の第一歩



- 9大学2研究所が共同運営し、全国共同利用に供する、データ科学、データ駆動科学、データ活用にフォーカスした高性能仮想化環境
- 2021年3月稼働 @ 東京大学柏IIキャンパス
- 現在ユーザ利用へ向けて詰めの作業を実施中



mdxハードウェアスペック

汎用CPUノード

Ice Lake 2 x 368ノード

総理論演算性能(倍精度) : 2.1PFLOPS

総メモリバンド幅 : 150.7TByte/秒

演算加速GPUノード

NVIDIA A100 x 8 x 40ノード

総理論演算性能(倍精度) : 6.4PFLOPS

総理論演算性能(半精度) : 100.7PFLOPS

総メモリバンド幅 : 496.3TByte/秒

高速NVMeストレージ

Lustreファイルシステム

1.0PB (NVMe SSD)

252GByte/秒

対外接続ルータ (Ethernet 100Gbps)

Ethernetネットワーク (100GbE/25GbE)

x368

x40

1.0PB

16.3PB

10.3PB

RDMA/ストレージ/仮想ディスクストレージネットワーク (100GbE, RDMA通信用)

大容量HDDストレージ

Lustre並列ファイルシステム

16.3PB

157.5GByte/秒

外部共有オブジェクトストレージ

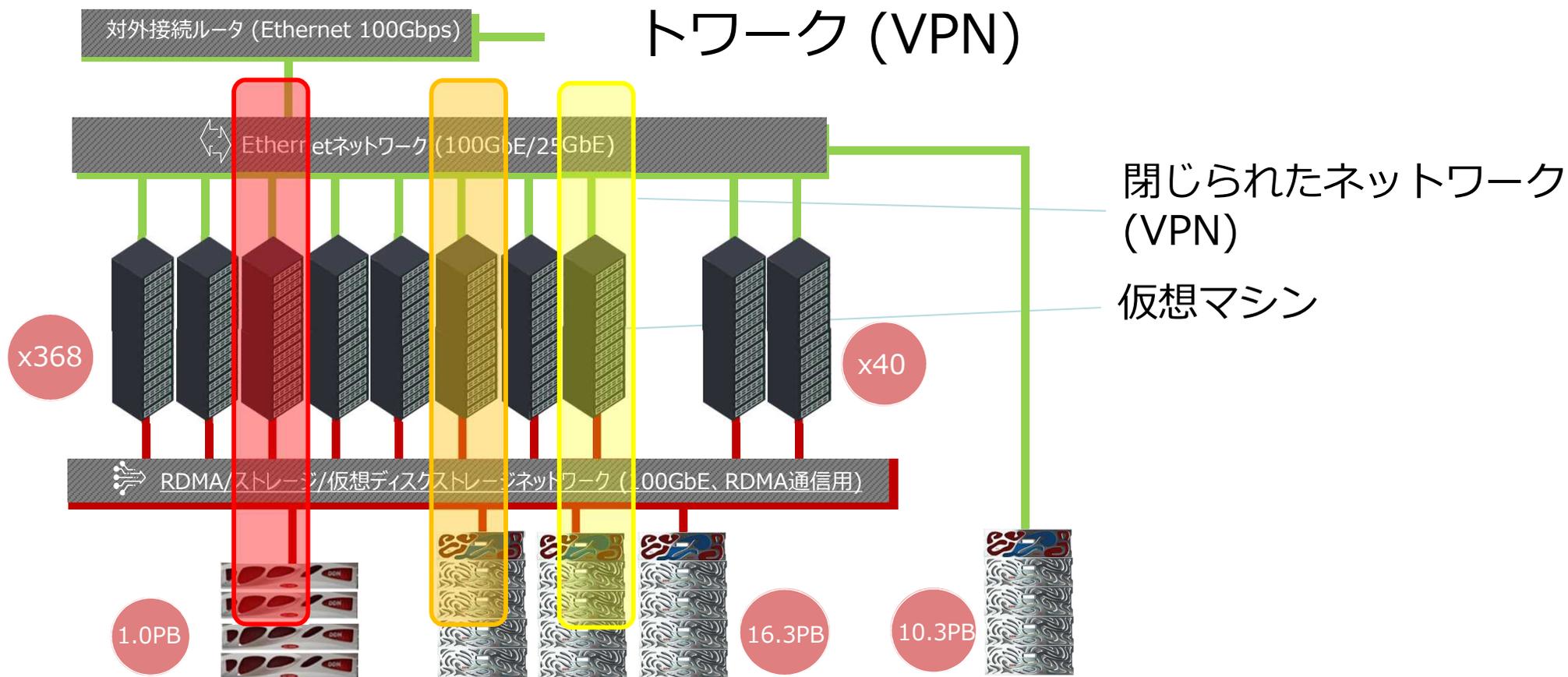
S3 Data Service

10.3PB

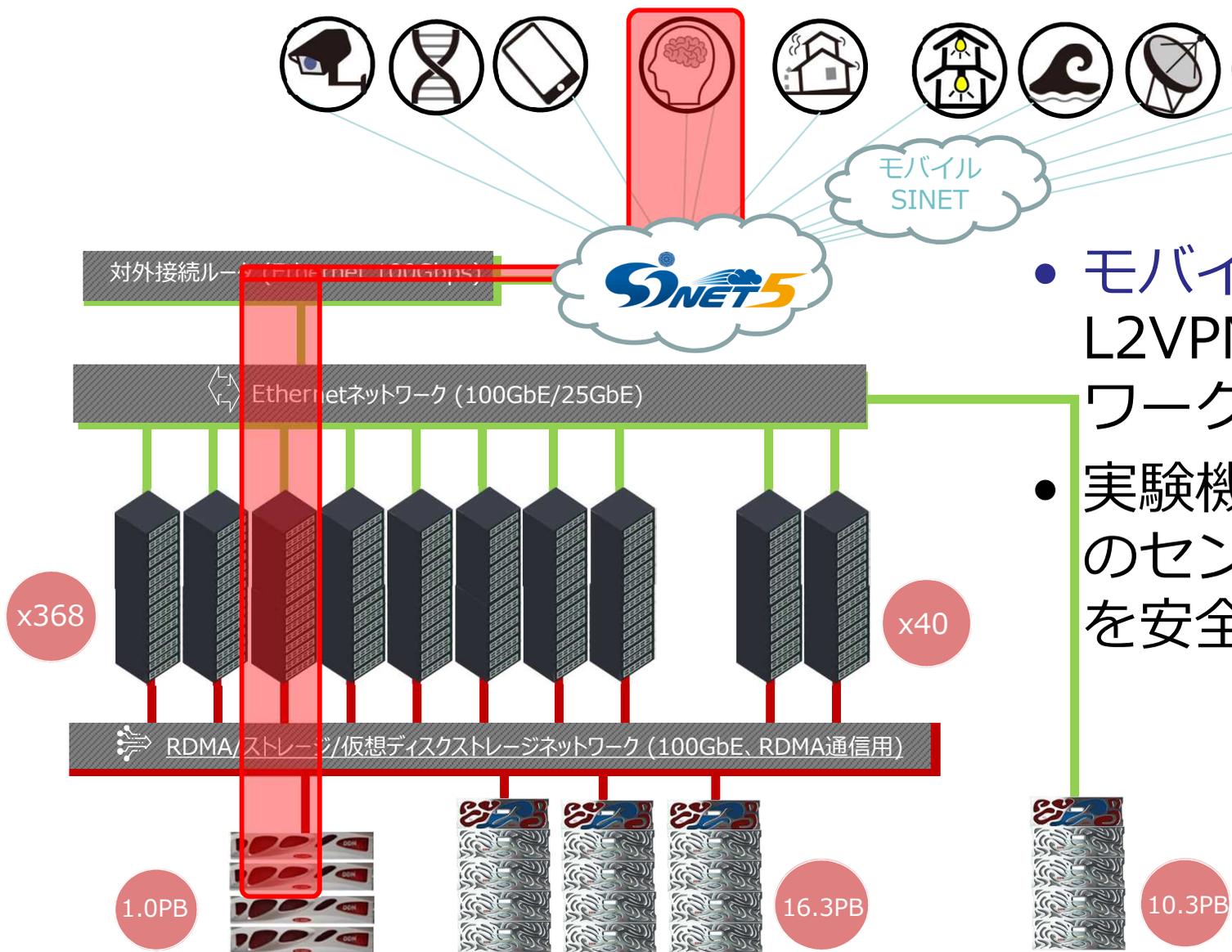
63.0 GByte/秒

仮想化環境

- 仮想化(VMware)により利用グループごとに占有環境(テナント)を提供
- テナント = 仮想マシン + 専用ネットワーク (VPN)



SINET・モバイルSINETとの連携



- モバイルSINET ≈ L2VPNをモバイルネットワークに延伸
- 実験機器, 野外や実験室のセンサなどからデータを安全に収集可能

ユーザビリティ: 様々な利用深度を想定

- 小規模・対話的利用 (≈ Kaggle, Colab)
 - 規定の環境 (Jupyterhubなど) で直ちに利用開始
 - Gakunin RDMで共有したデータを処理する環境をBinderHubで起動
 - パブリックサービスにない自由度 (SSHなど) も提供
- 独自環境 (≈ AWS IaaS)
 - 独自にカスタマイズした環境 (Jupyter, 特定MLフレームワーク、etc.)
 - その他、パブリックなIaaS同様、自由に環境を構築可能
 - そのための仮想マシンのテンプレート、そのレシピを公開 (コミュニティで開発・共有)
- 中規模クラスタ環境
 - 多数のノードからなる並列・分散処理環境 (ビッグデータ、中規模AI・ML、データ収集・蓄積・サービス環境)
 - そのためのレシピも共有

ユーザビリティ: 柔軟性・自由度

- クラウド風の使い方(バッチスケジューラではなく)
- ≈ 仮想マシンを長時間, 常時稼働可能
 - = 申請した分の仮想マシンは一定時間以内に起動すること(実効的な常時稼働)を保証 ⇒ 起動保証VM
- + 「資源が空いている間のみ稼働できる」仮想マシン ⇒ スポットVM
 - 計算資源が必要, 実時間性は不要なジョブ(≈ バッチジョブ)
 - 連携マシン(ABCI, Wisteria/BDEC01)での実行も可能
- 外部通信
 - 各テナントが外部(インターネット)との通信を制御可能

つなげるを前提とした設計

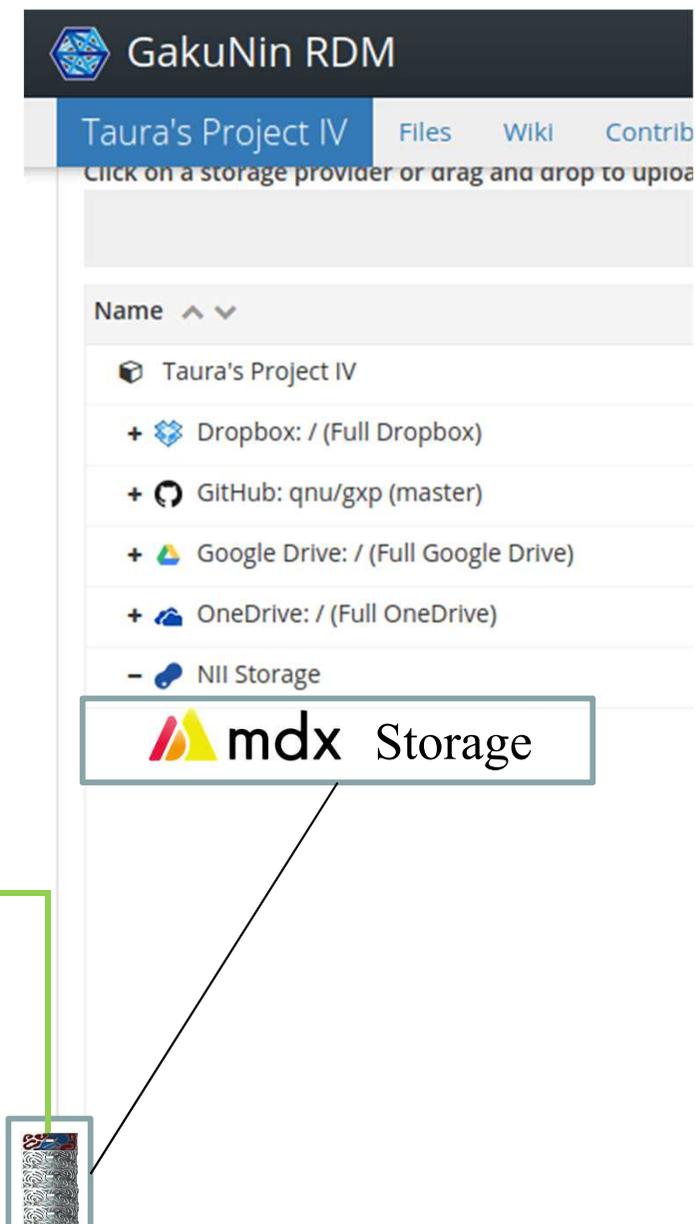
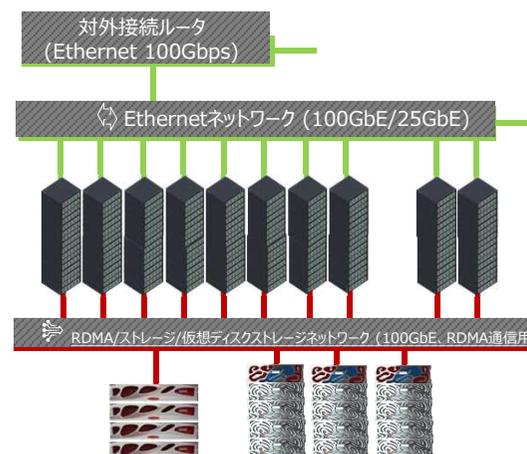
-  GakuNin 国内多くの大学・研究機関ですぐにサインアップ可能
-  GakuNin RDM と連携
 - 様々なストレージ, データ・コードレポジトリとの連携機能
-  mdx
 - Gakunin RDMの外部ストレージとして連携
 - まずはS3共有オブジェクトストレージ
 - (次) 内部大容量ストレージ(Unixファイルシステム)との連携
 - (進行中) BinderHubを用いたストレージ-計算機の連携

公開基盤(JAIRO), 検索基盤
(CiNii)との連携を協議中

Gakunin RDMを用いたデータ共有とmdx

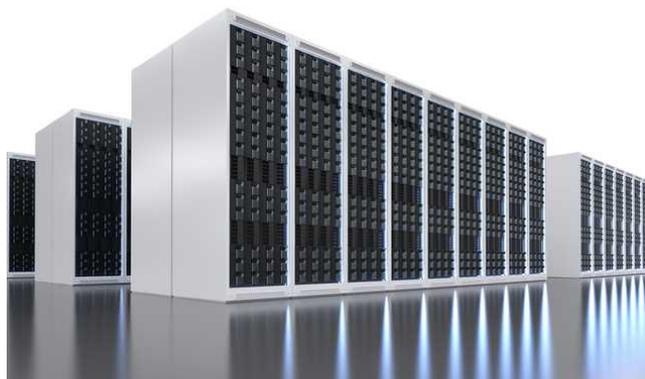
1. 「プロジェクト」を作成
2. プロジェクトにストレージを接続
 - Google Drive, Microsoft OneDrive, Dropbox, Githubなどを接続可能
 - Nextcloud, S3互換ストレージなど、オンプレスト
レージの接続も可能
 - NIIデフォルトストレージを提供
3. プロジェクトに共同研究者を招待
 - 参考: 北大、名古屋大
<https://www.youtube.com/watch?v=SzS8-o5B3vw>

mdxがGakunin RDMと接続するストレージと、それを処理する計算基盤を提供



将来像 NIIと情報基盤センター群のマシン間データ連携

- mdxは端緒に過ぎない
- 多数の大学の共同利用マシンが Gakunin RDMと接続、互いにも連携した状態を目指す
 - 大規模データの共有～直接処理までを迅速にできる広域クラウド環境



GakuNin RDM

Taura's Project IV Files Wiki Contrib

Click on a storage provider or drag and drop to upload

Name ^ v

- + Taura's Project IV
- + Dropbox: / (Full Dropbox)
- + GitHub: qnu/gxp (master)
- + Google Drive: / (Full Google Drive)
- + OneDrive: / (Full OneDrive)
- NII Storage

mdx Storage

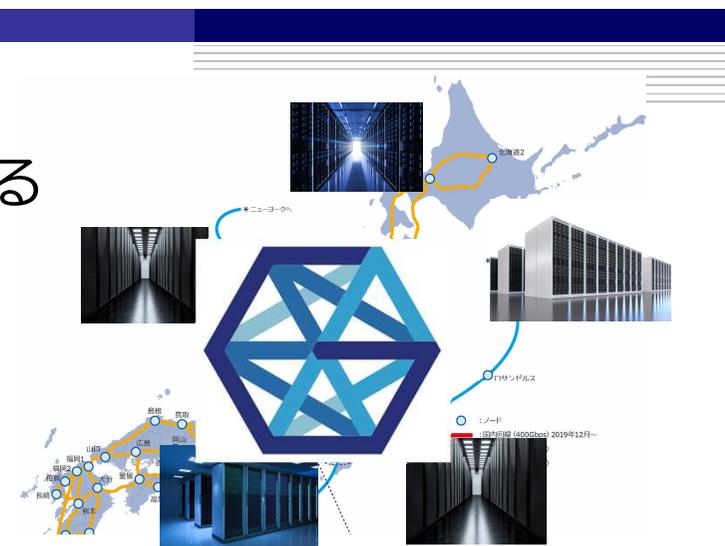
- 大学 □□
- △△大学 ●●
- ◎◎大学 ◆◆
- ▲▲大学 ☆☆
- ...

mdx導入経験の共有・継承

- 総合評価入札では（きっと）できなかった
 - 総合評価入札の限界
 - 「製品化」されたものしか提案できない
 - 開札から納入までの期間の短さ
 - お互い（大学、業者とも）「経験済みの領域」でしか導入できない
 - 学術基盤を進化させていく上で検討すべき課題
- 多機関連携（9大学2研究所）で設計の議論～調達～運営までの経験をすべて共有

まとめに変えて: 次世代学術基盤に向け, データ基盤, 情報基盤センター群, NIIの果たす役割

- 学認など全国的インフラを利用した迅速感
- ソフトウェアの急速な進化 (変化?) に追従できる柔軟な環境
- データ隔離・セキュリティ
- 利用をハブとしたコミュニティ創成



活用分野

共通ディシプリン: 計算科学

原理のシミュレーション中心
核物理、材料・物性、生体分子、
燃焼、気象、地震、宇宙、etc.

共通ディシプリン: データ科学

データ解析、モデリング中心
医療、経済、空間、モビリティ、
環境、歴史 (アーカイブ等)、
ソーシャルデータ、etc.

情報基盤

高性能計算機、ストレージ (スパコン) 中心

データ共有基盤、データ解析環境、データプラットフォーム構築、実時間データ収集、モバイルネットワーク

情報系研究分野

数値計算アルゴリズム、HPC中心

機械学習、NLP, 画像処理、大規模データ処理、ネットワーク、クラウド、etc.