

# mdx : データ活用のためのプラットフォームと, 医療データでの活用について

---

田浦健次郎

東京大学 情報基盤センター長

大学院情報理工学系研究科 教授

JHPCN共同利用共同研究拠点 拠点長

[tau@eidos.ic.i.u-tokyo.ac.jp](mailto:tau@eidos.ic.i.u-tokyo.ac.jp)



# 謝辞

- mdx 構想・運営に参画している以下の機関の皆様



- 日々運用、改善に尽力されている皆様
- 利用者の皆様
- すべての関係者の皆様に感謝いたします

# 自己紹介

- 東京大学
  - 情報理工学系研究科
  - 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
- 研究分野：システムソフトウェア（並列処理、高性能計算、並列データ処理系、（最近）プライバシー保護）

# 本日の目次

## ▪ mdx

- ≈ より広範な分野に使える・使いやすい、クラウド型の計算機
- これまでほぼ大規模高性能シミュレーション/深層学習のためのスパコン ⇒ データ利活用、データプラットフォーム構築、Society 5.0 実現へ向けたサービス構築

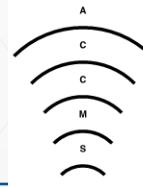
## ▪ 共同利用共同研究拠点JHPCN

- データ科学・データ活用を軸とした情報学 - 様々な分野との学際連携
- これまほほ計算科学 ⇒ {計算科学+データ科学・データ活用}

## ▪ プライバシー保護のためのシステムソフトウェア研究

# mdx

- 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・データ活用応用にフォーカスした高性能仮想化環境 <https://mdx.jp/>
- @ 東京大学柏IIキャンパス



# mdx

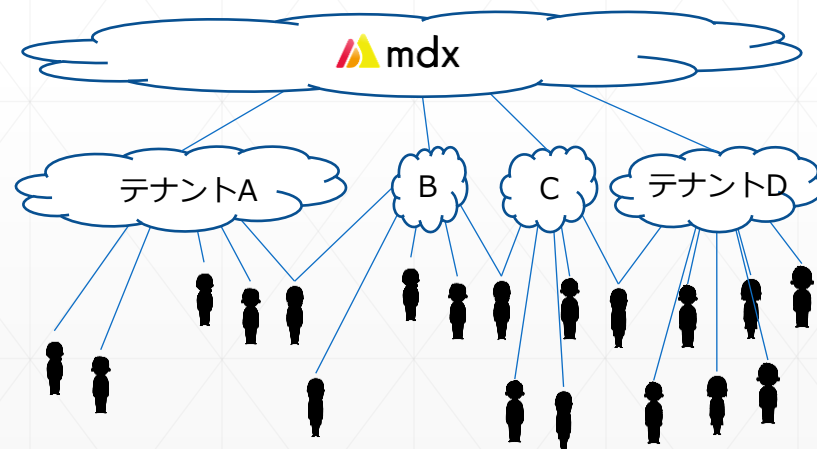
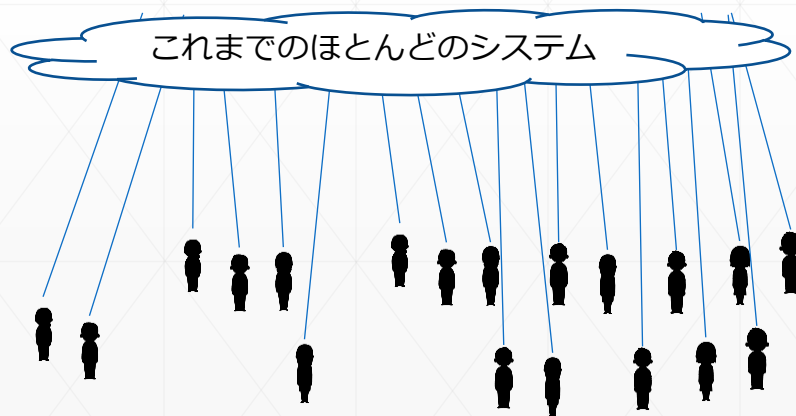
- 申請ポータルから利用申請（学認でサインイン）
- 申請が受理されたらユーザポータルから「仮想マシン」を生成
- 自分（たち）だけの環境が使える ⇒
  - 自前で調達したサーバと同様、自由に環境を構築可能
  - 他のユーザ環境と隔離されている（高セキュリティ）
- ハードウェア的にはそれなりの規模のものを共有 ⇒
  - 簡単・迅速に仮想マシンを追加（環境の拡張）可能
- 「Infrastructure as a Service」型のクラウド環境
- 学認でサインインし、ほかのサービス（e.g., Gakunin RDM）とシングルサインオンで連携

# 学認でサインインできます、が…

- 日本の多くの教育・研究機関の大学アカウントが学認と連携しています
- したがって多くの教育・研究機関の方がすぐに申請ページにアクセスでき（るようになり）ます
- しかし、各機関で学認連携担当者が「mdxというサービスの利用許可」を出す必要があります
- もしmdxにアクセスして「このサービスにはアクセスできません」的なエラーになったら、学認連携部署に「使いたい!」とご連絡いただくと幸いです
- 設定に必要な情報：「学認 SP 一覧」などで検索  
<https://www.gakunin.jp/participants>

# これまでのマシン（スパコン）とmdxの違い

- これまで（≈ いわゆるスーパーコンピュータ）
  - 全員が、管理者が設定する**単一の環境**を使う
  - 柔軟性がなく、目的、したがってユーザの分野が限定
- mdx：**仮想化された高性能環境**
  - 分野、グループごとに異なる環境を、それぞれのユーザグループで構築可能

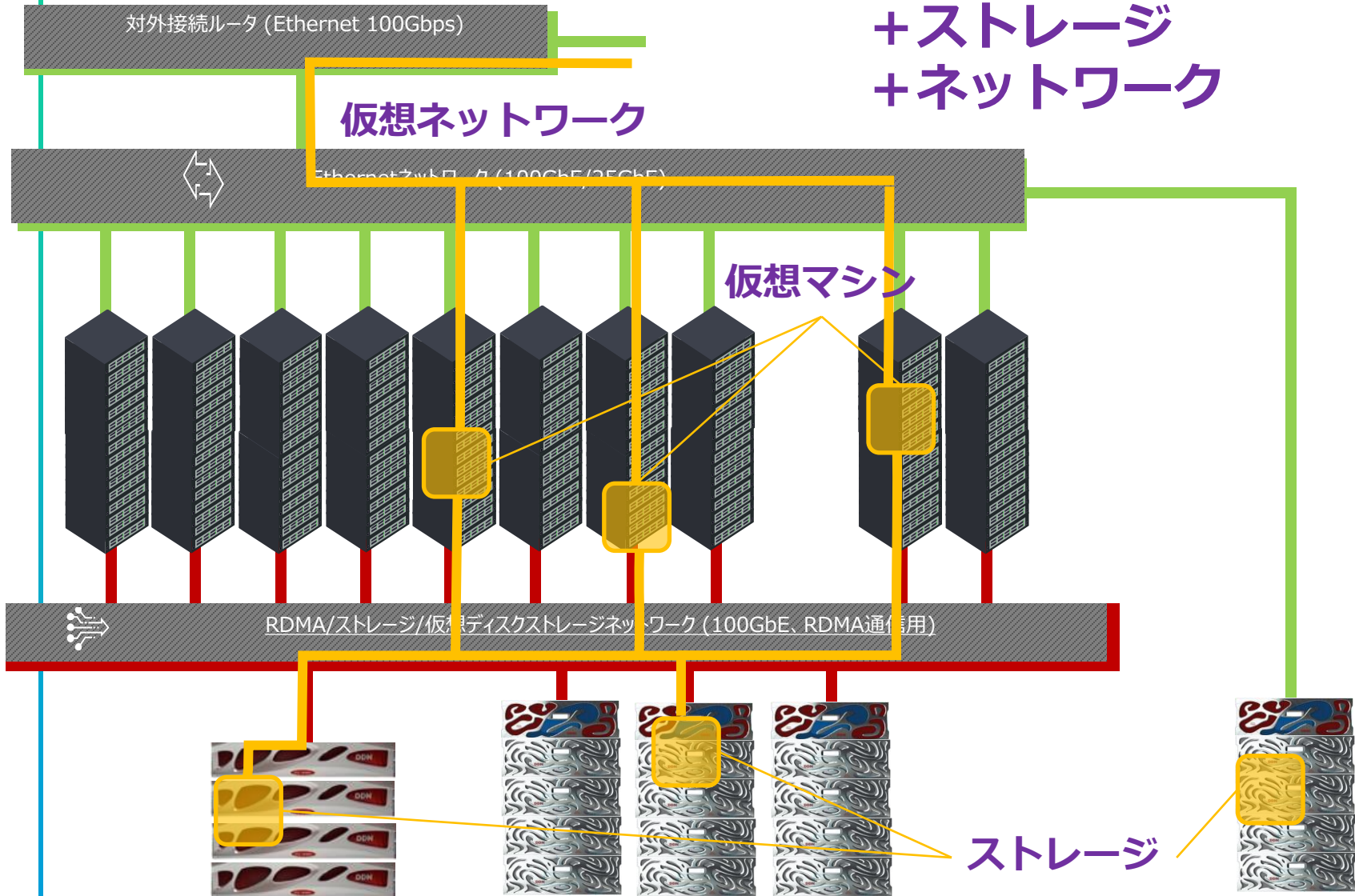


スパコン用の超高性能計算は無縁な方にも有用です

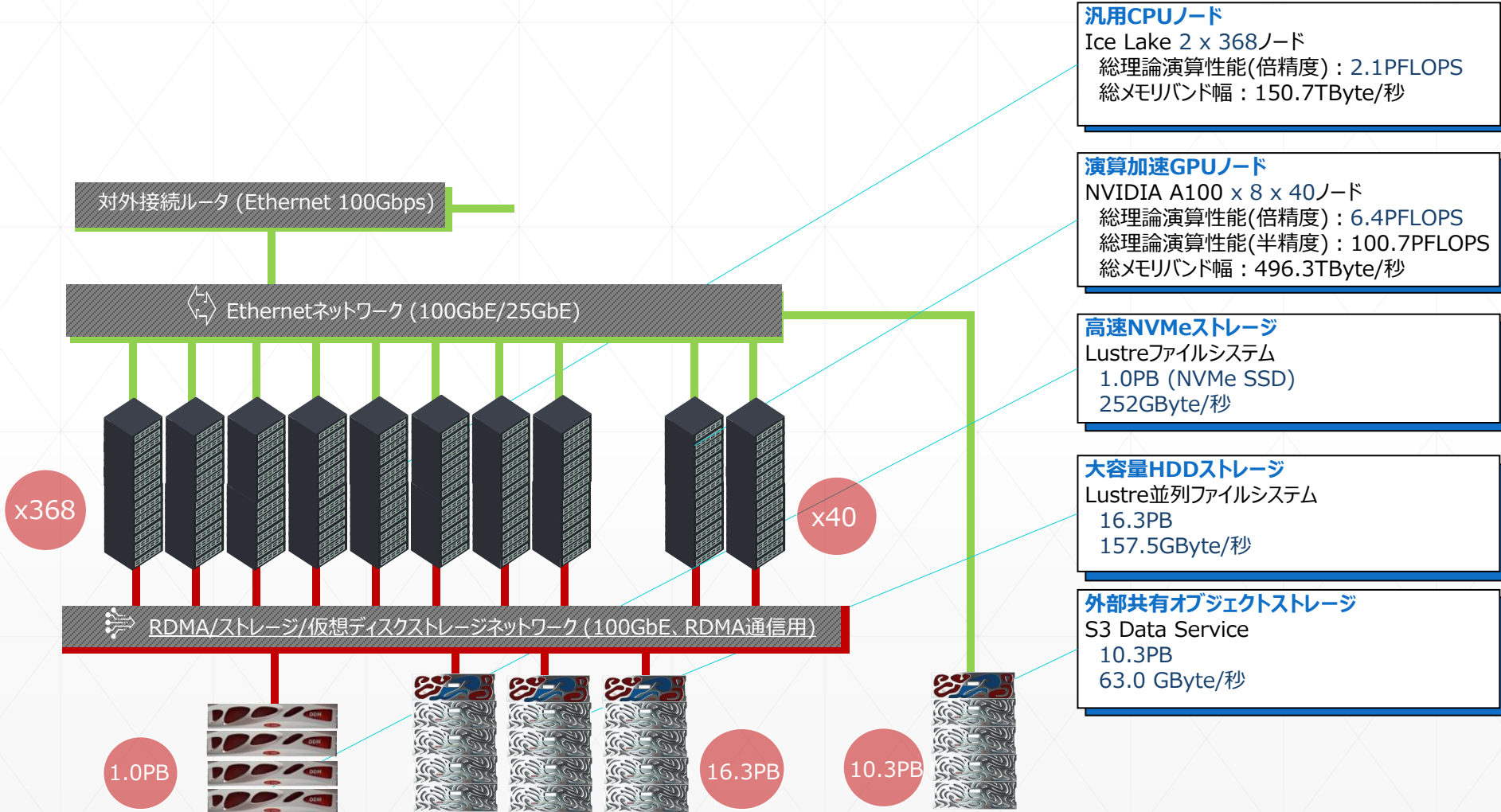


# mdx利用モデル (図解)

テナント  
≈ 仮想マシン  
+ ストレージ  
+ ネットワーク

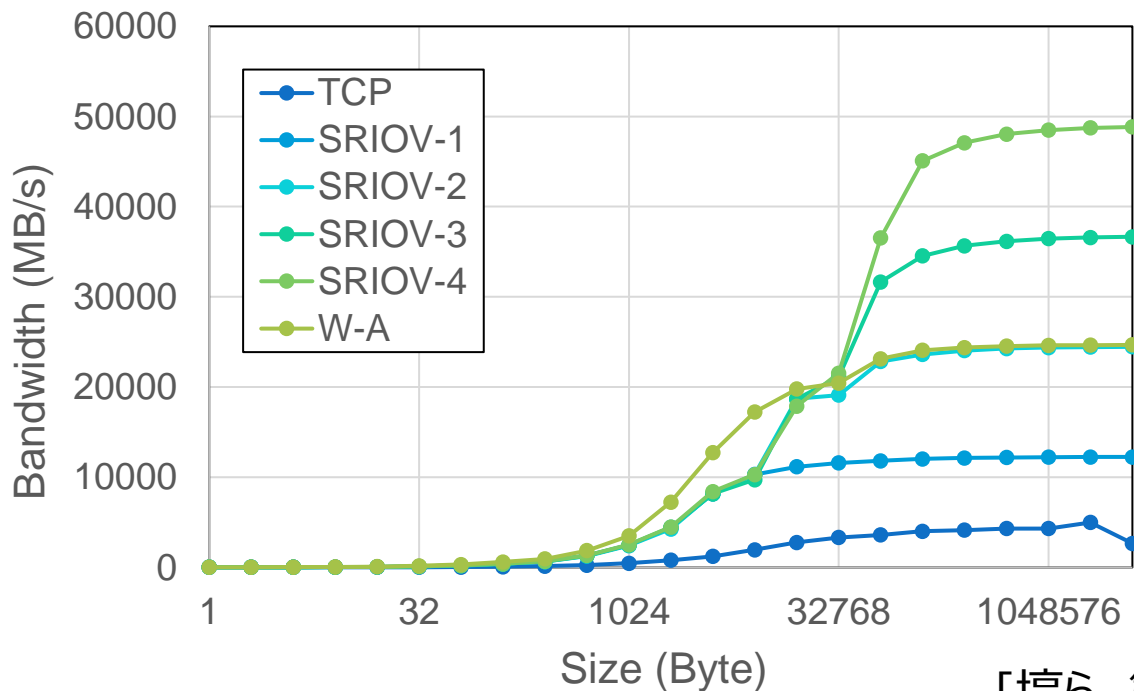


# mdxハードウェアスペック



# 通信性能

- 100Gbps RDMA over Converged Ethernet
- GPUノードには 4つのNIC (合計 400Gbps)



- 仮想環境であるが通信性能も高く、中規模までの並列計算にもスケール可能と期待

[埴ら、第183回HPC研究会]

# 仮想化 ⇒ 使いやすさ・セキュリティ向上

## 1. 仮想化による柔軟性

- ホストを（管理者権限で）自由に設定可能

## 2. VPN (VLAN) による柔軟性

- プロジェクトごとにネットワークの延伸が可能

## 3. 仮想化・VPNによるセキュリティ（分離）

- 他のグループの環境のオペミスや脆弱性が他に影響しない

## 4. 環境の互換性・可搬性を高められる

- VMテンプレート (cf. MateriApps)

# 環境自動構築サポート (Software Defined Environment)

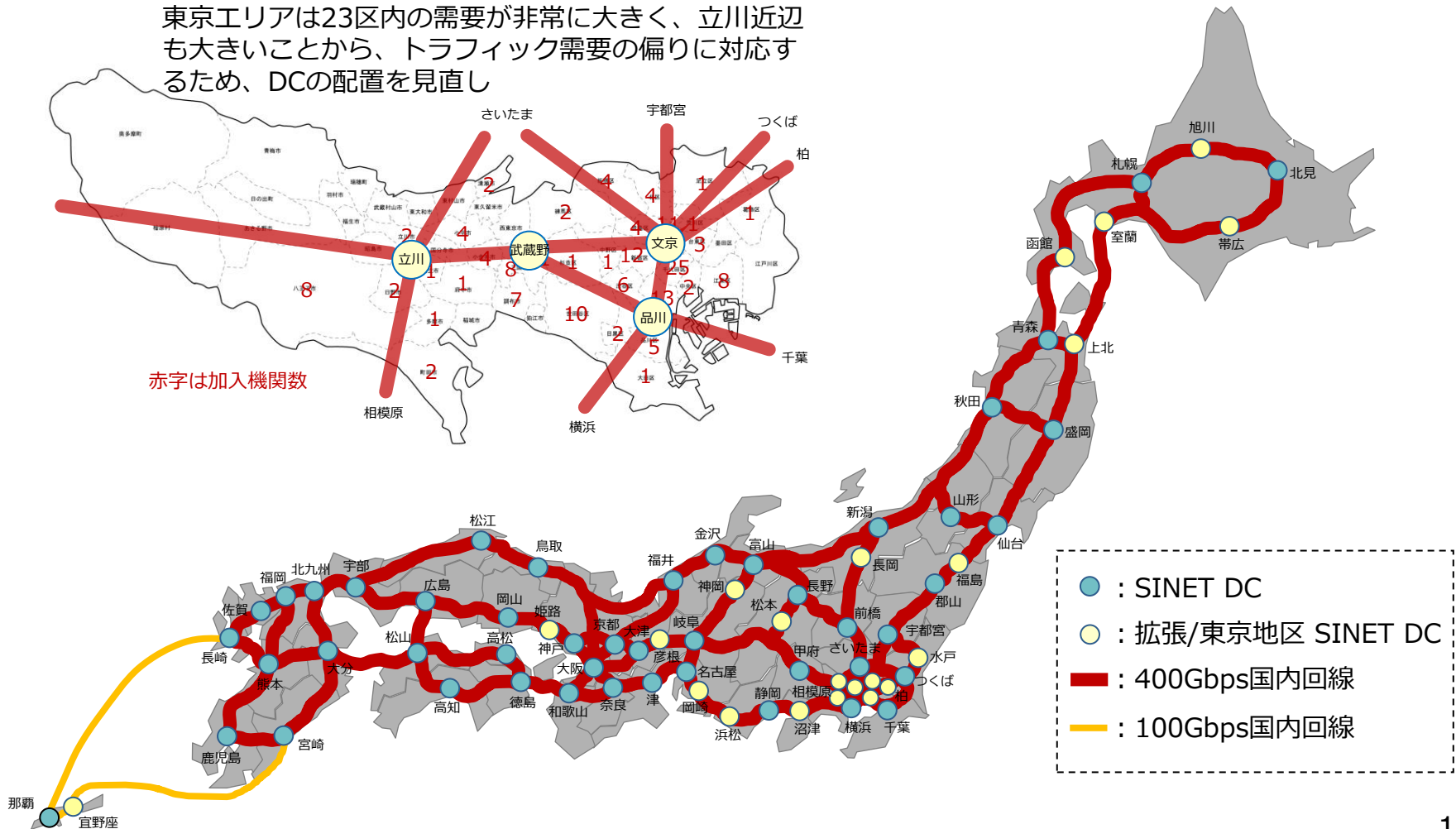
- mdx用Ansible playbook
  - Special thanks to: 中村遼先生@東大、空閑洋平先生@東大、杉木章義先生@北大
- Kubernetes環境
  - 杉木章義先生@北大
- Jupyter (PaaS) 環境
  - 華井雅俊先生@東大
- mdx REST API
  - 大江和一先生@NII, 竹房あつ子先生@NII, 工藤知宏先生@東大

オンデマンドに環境を伸縮  
環境の移植性  
Coming soon

- 国立情報学研究所が運用する、研究機関を接続する学術ネットワーク
- 国内ほぼ全県を400Gbps\*で接続

\* 沖縄は技術的な制約により当面100Gbpsベース

東京エリアは23区内の需要が非常に大きく、立川近辺も大きいことから、トラフィック需要の偏りに対応するため、DCの配置を見直し



# SINET VPN

図提供：国立情報学研究所

- インターネットに加え閉域網が利用可能で、共考共創の高度なサービスも活用可能

**大型実験施設等**  
(全国各地、海外各地)

図は例

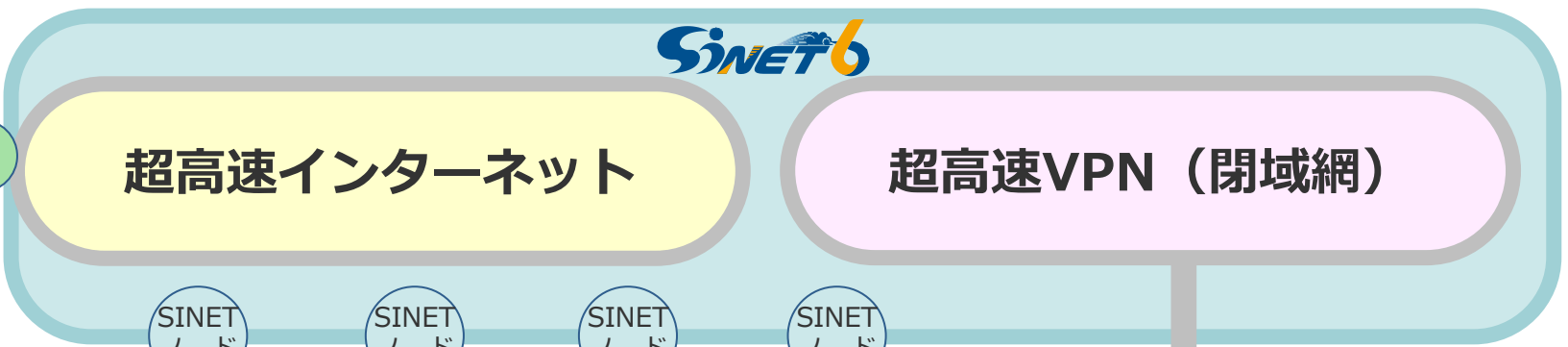
**スパコン**  
(HPCI 13拠点)

**mdx**

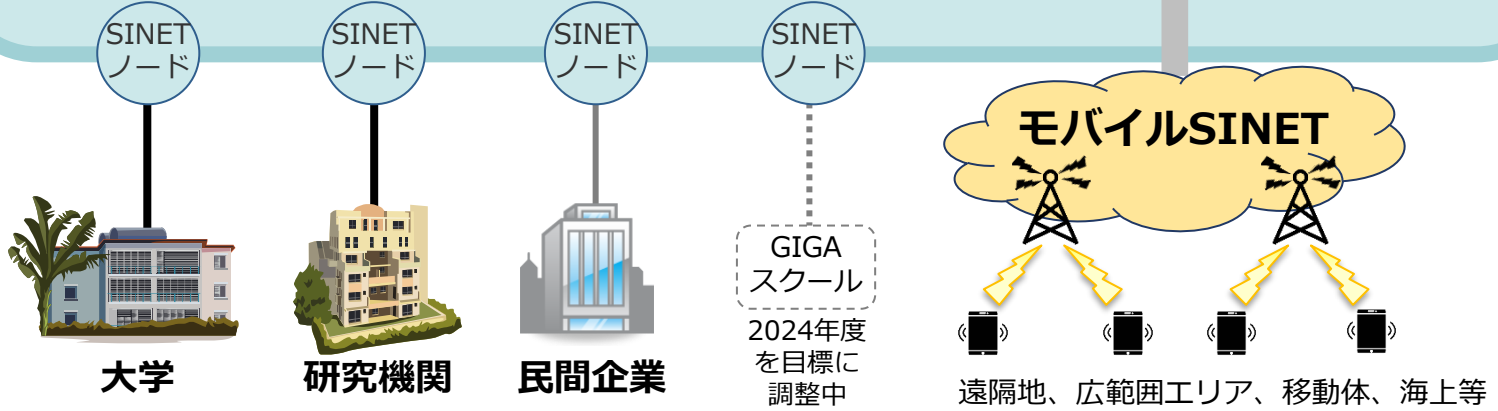
**研究データ  
基盤**

**直結クラウド**  
(34社45拠点)

- Web会議システム
- クラウド
- 商用ISP
- 商用ISP

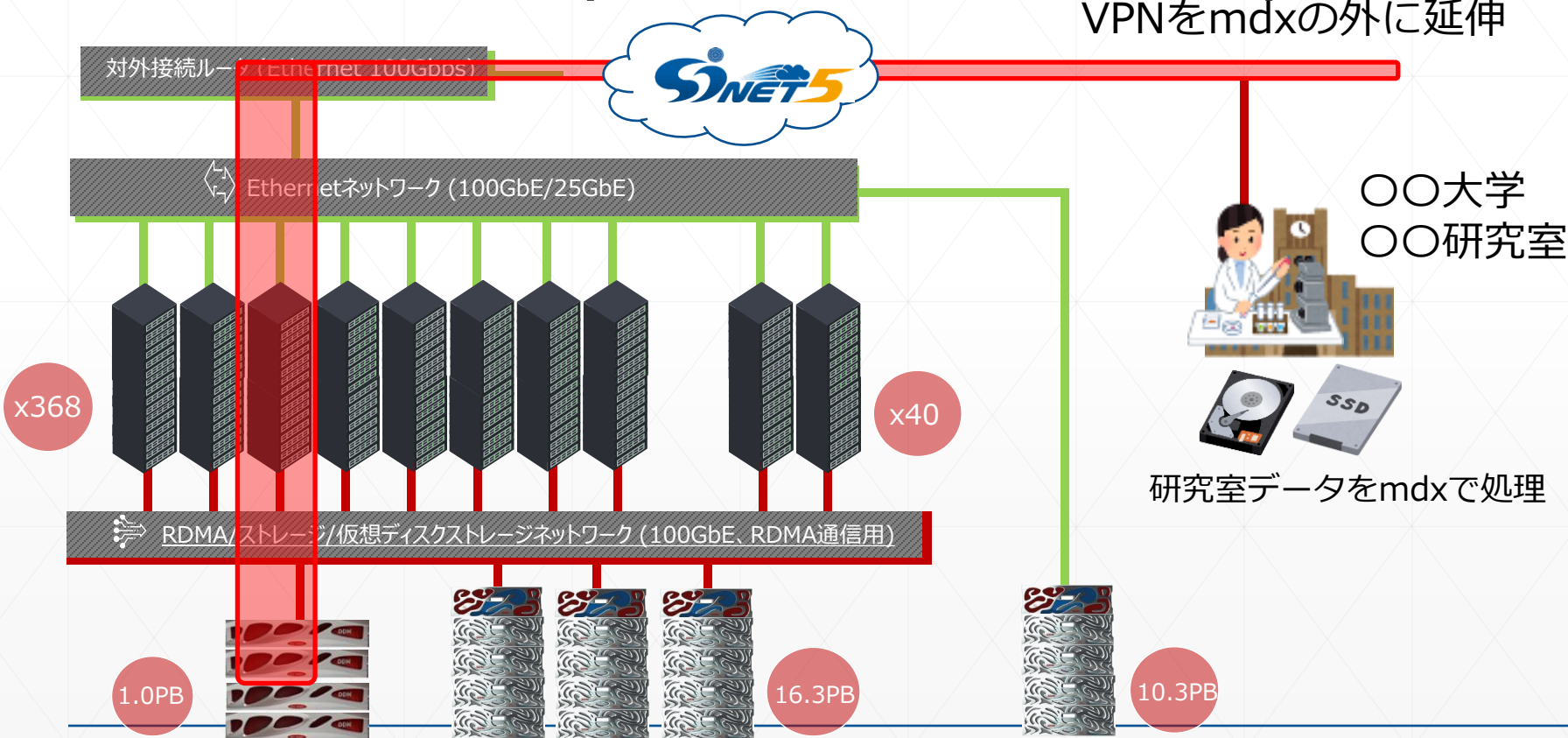


**990 機関**  
(2022年3月末現在)



# SINET VPNとmdx

- SINET VPNによりmdx内に構築した環境を研究室, 他大学のマシンと(同一LANにいるかのように)接続可能  
VPNをmdxの外に延伸

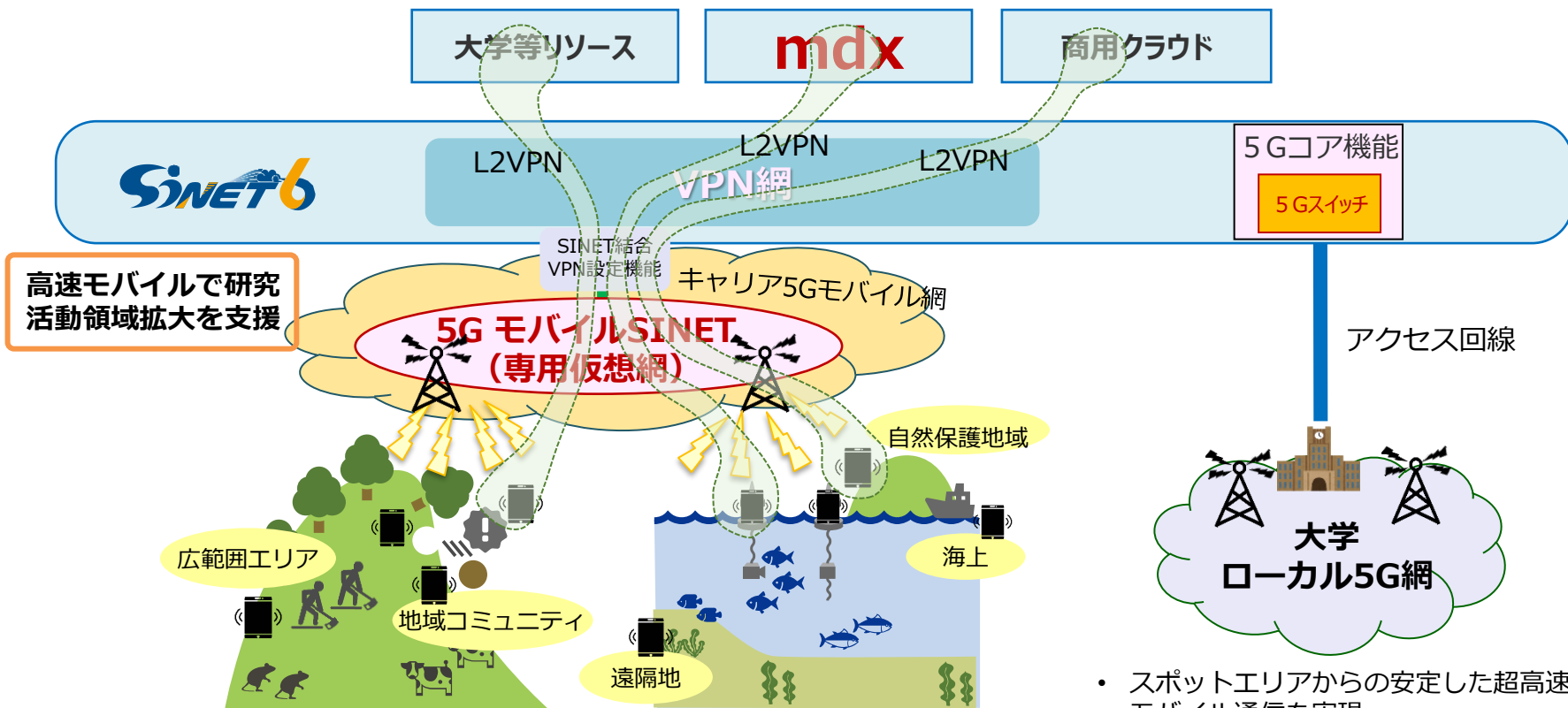




# mdxとモバイルSINET

図提供：国立情報学研究所

- モバイルSINET：商用モバイル網の中にSINET専用の仮想5G網（4G/3Gも可）を形成してSINET VPN網と接続することで、セキュアな通信環境を実現
- ローカル5G：SINET側に5Gコア機能を実装することで大学のローカル5G網の経済的な導入を支援（まずは小さな規模で実証実験を開始予定）

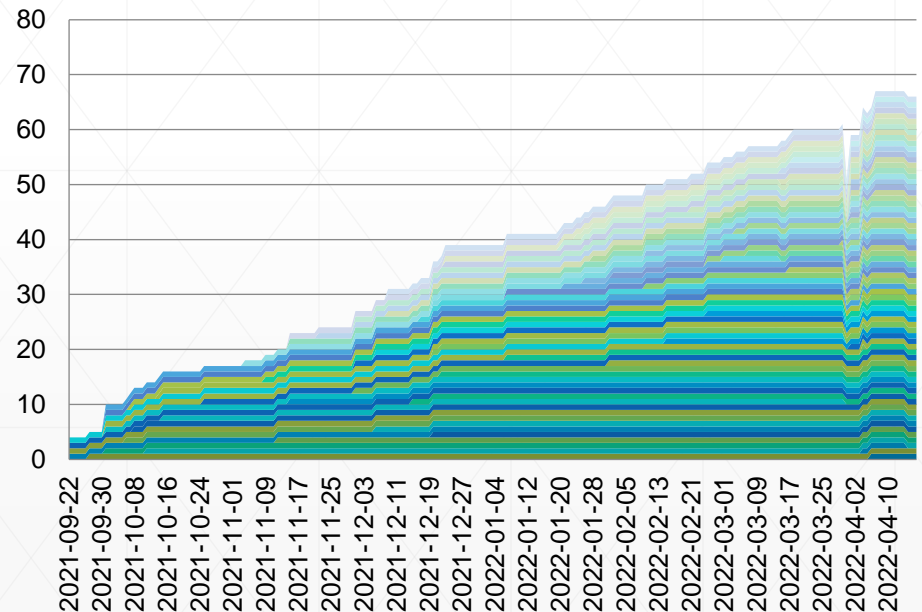


- 遠隔地、広範囲エリア、移動体、海上等を広くカバー
- セキュアな通信環境を実現

- スポットエリアからの安定した超高速モバイル通信を実現

## ユーザグループ (プロジェクト) 数

- 深層学習を応用したX線レーザーイメージング画像のデノイジング
- 建物の築年と構造の推定モデルの開発
- mdxマテリアルリサーチ連携
- 労働経済学プロジェクト
- 疑似人流データの開発
- DIAS連携
- 合成人口プロジェクト
- ...



# 共同利用共同研究拠点 JHPCN

# 共同利用・共同研究拠点の枠組み

- [https://www.mext.go.jp/a\\_menu/kyoten/](https://www.mext.go.jp/a_menu/kyoten/)
  - 「個々の大学の枠を越えて大型の研究設備や大量の資料・データ等を全国の研究者が共同で利用したり、共同研究を行う」仕組み
- 現在100拠点が認定されている（一部は国際共同利用・共同研究拠点）
  - 物理系：
    - 宇宙線研（東大）、理論物理学研究拠点（京大）、…
  - 情報系：
    - JHPCN（8大学）、計算科学研究センター（筑波大学）、…
  - …

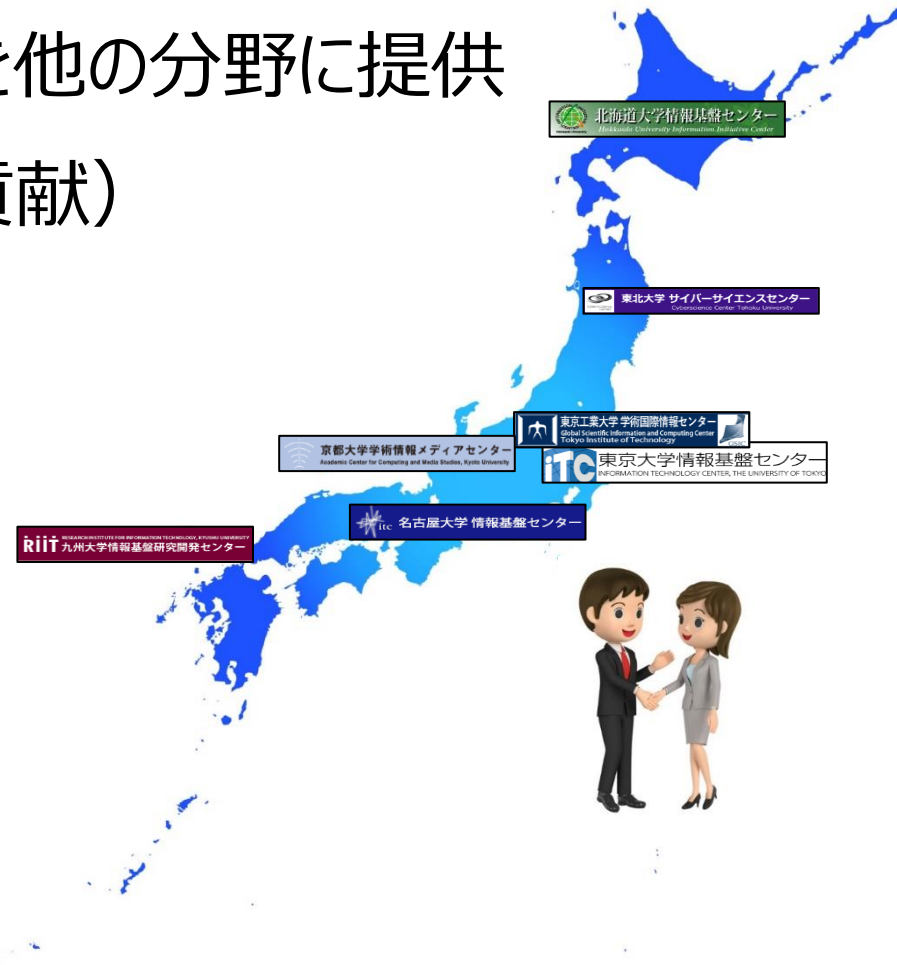
# JHPCN（正式名：学際大規模情報基盤 共同利用共同研究拠点）

- 8大学の基盤センターが運営する共同研究のための組織
  - 北大、東北大、東大、東工大、名大、京大、阪大、九大
- 全国から共同研究を募集
  - 現在高性能計算シミュレーション中心の分野が多く集まる
  - 採択された課題に計算機利用時間割り当て
- 2022年度より「データ科学・利活用分野」募集を開始
- <https://jhpcn-kyoten.itc.u-tokyo.ac.jp/>

# 基盤センターとJHPCN拠点の使命

- 情報の専門家と情報基盤を他の分野に提供
- 情報学 x ○○学（学際貢献）
- コミュニティ形成

mdxはこの使命をこれまでのシミュレーション中心・計算科学から、**分野も基盤もセクターも**「データ駆動科学・データ活用分野」に広げるための第一歩



(注: リンク先ページはmdx利用以外も含む. 以下はmdx利用)

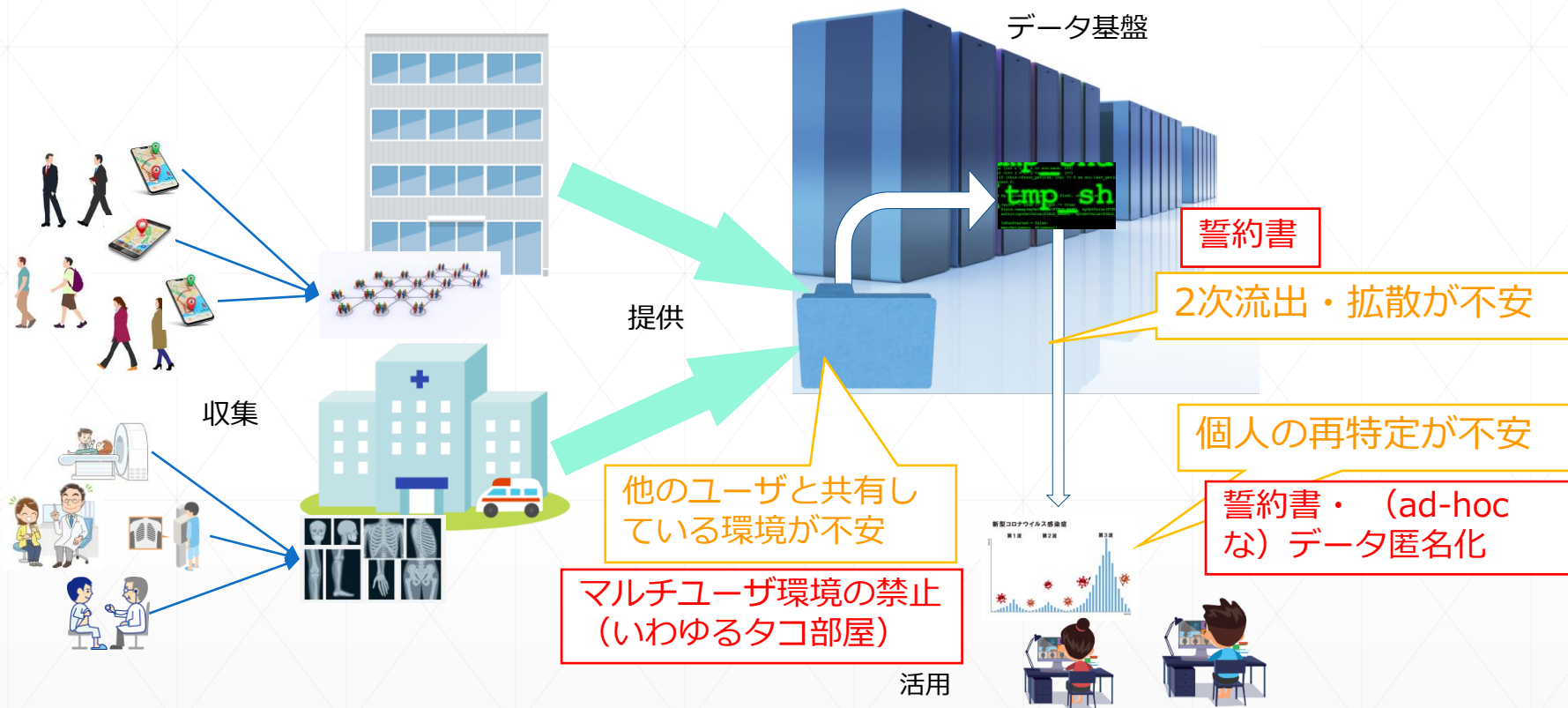
## JHPCNでのmdx利用課題

- 財務ビッグデータの可視化と統計モデリング
- 医療・介護領域の人材マッチングに最適化された大規模グラフニューラルネットワーク
- エージェントモデルと統計データを用いた全国規模の疑似人流データの開発
- 大規模な日本語モデル構築・共有のためのプラットフォームの形成
- グラフニューラルネットワークとマルチタスク学習による汎用的物性予測モデルの構築
- ビヨンド・"ゼロカーボン"を目指し地域と技術をつなぐ情報基盤の構築
- 単語間に区切りのない書写言語における係り受け解析エンジンの開発
- 多次元高精細地表情報 (MHESD) の地球科学・歴史考古学における高度利活用

# データセキュリティ・プライバシー 保護データ解析研究



# データ提供・利活用のボトルネック



データ提供が小規模にとどまる（萎縮の）背景

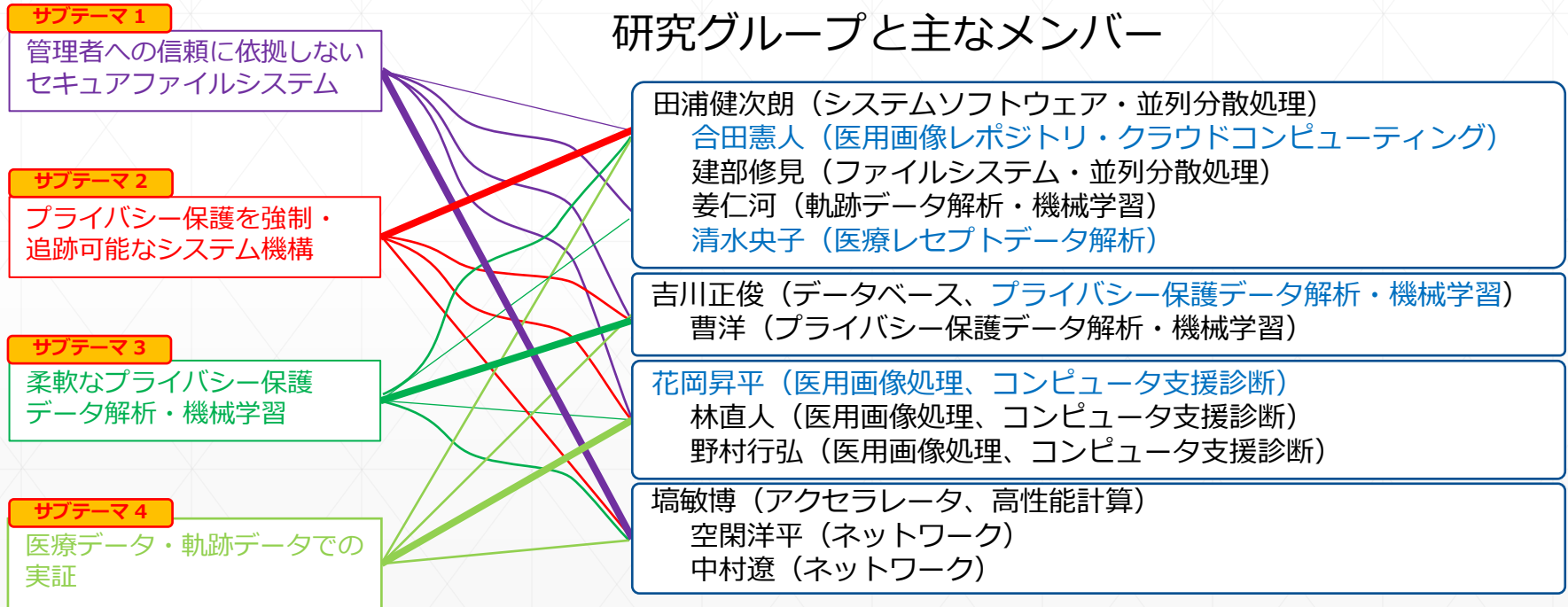
## 本研究が目指す基盤技術

- データ提供者の懸念を解消し、利活用を促進するIT基盤技術
- JST CREST 「基礎理論とシステム基盤技術の融合によるSociety 5.0のための基盤ソフトウェアの創出」領域「実応用に即したプライバシー保護解析とセキュアデータ基盤」として実施

# 研究体制

- 基盤ソフトウェア、システム運用、データベース・プライバシー、軌跡データ解析、医療データ解析の研究実績のあるチームで一体的に推進する

## 研究グループと主なメンバー

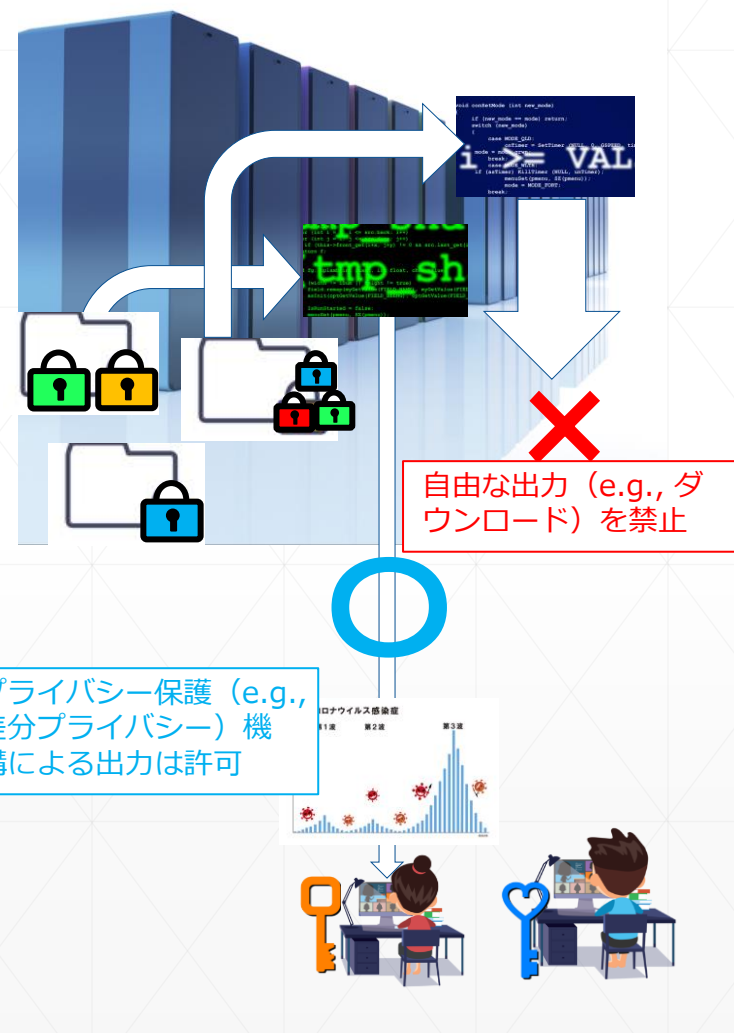


# データ提供者のジレンマ

- データを汚す
  - ad-hocなデータの匿名化 ⇒ 不安が消えない、突き詰めるほどデータの有用性が損なわれる
  - ローカル差分プライバシー ⇒ データの有用性が損なわれる
- 結果を（確率的に）汚す
  - グローバル差分プライバシー ⇒ 解析者を信用する必要がある
- ⇒ 解析結果がプライバシー侵害していない（適切な差分プライバシー機構を通して）ことをデータ解析者への信用に委ねずに行える機構（本研究のメインテーマ）

# 本研究が目指す基盤技術

- データ提供者が安心して解析プラットフォームにデータを提供できるようにするための基盤技術
  - 管理者への信頼に依拠しないアクセス制御が可能なファイルシステム（暗号化ファイルシステム）
  - データ提供者は後から共有を取り消し可能
  - 解析プラットフォームからデータの持ち出しが不可能・プライバシー侵害を起こさない出力だけを持ち出し可能（サンドボックス）



# プライバシー機構の設計方針

- プログラムからのデータの出力を制限
- ただし、出力してよいかどうかをデータだけを見てシステムソフトウェアに賢く判定させることは無理、非現実的
  - 数値3.54を見て「ノイズが載っているか」はわからない
  - ましてやどれだけ載っているかなどわかるはずない
  - 載っている量がプライバシーを侵害しない適切な量などもデータだけで判断できるものではない

- ⇒ 実際は「勝手な出力を禁止、特定の（信頼（できれば検証）された）プライバシー機構を通した出力のみ許可できる」枠組み

# まとめ

- mdxは柔軟性、セキュリティを備え、「データプラットフォーム構築事業」の共通基盤を目指している
  - SINET VPNと連携して他のシステムと高速に、安全に接続も可能
- データ利活用、安全性、プライバシー保護の観点から「データ＋解析プラットフォーム」の提供というモデル
  - 解析（クエリの実行系）がプライバシー侵害を防ぐ
  - 解析プラットフォーム外へのデータの持ち出しを防ぐ

ご清聴ありがとうございました