

Tomo-e Gozen と mdx

Satoshi TAKITA (IoA/UT)

mdx とは

9 大学、2 研究機関が共同運営する
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

実体は東大、柏 II キャンパス

2021-09-22 より試験運用開始 (無償利用)

2023-05-10 より本運用 (課金開始)

「データ活用社会創成プラットフォーム」は
用途に応じてオンデマンドで短時間に構築・拡張・融合できる
データ収集・集積・解析機能を提供するプラットフォーム。

データ活用社会創成プラットフォーム 3本柱

- 1 SINETを活かしたリアルタイム収集・集積・解析環境の動的な構築**
遠隔地のセンサーやストレージ、データプラットフォームの計算資源、ストレージをつないで、リアルタイムに入力から出力を得られるアプリケーションごとの収集・集積・解析環境（仮想データプラットフォーム：仮想DP）を、使いたいときに即時に構築する
SINETモバイル基盤によりセンサー等のデータを安定してセキュアにつなぐ
- 2 高性能計算環境によるデータ科学と計算科学の融合**
データ科学、計算科学の手法を融合し、さらに国内最高の計算環境を用いて他に無い高精度の予測を行えるようにする
- 3 異種データ・異種知識の融合活用の推進と利用者支援**
様々な分野のデータ保持者、解析者、利用者が産学にまたがって連携するコミュニティを形成し、新たな価値創造につなげる。
データ活用を目指す利用者へのコンサルティングや開発支援を実施する。

<https://mdx.jp>

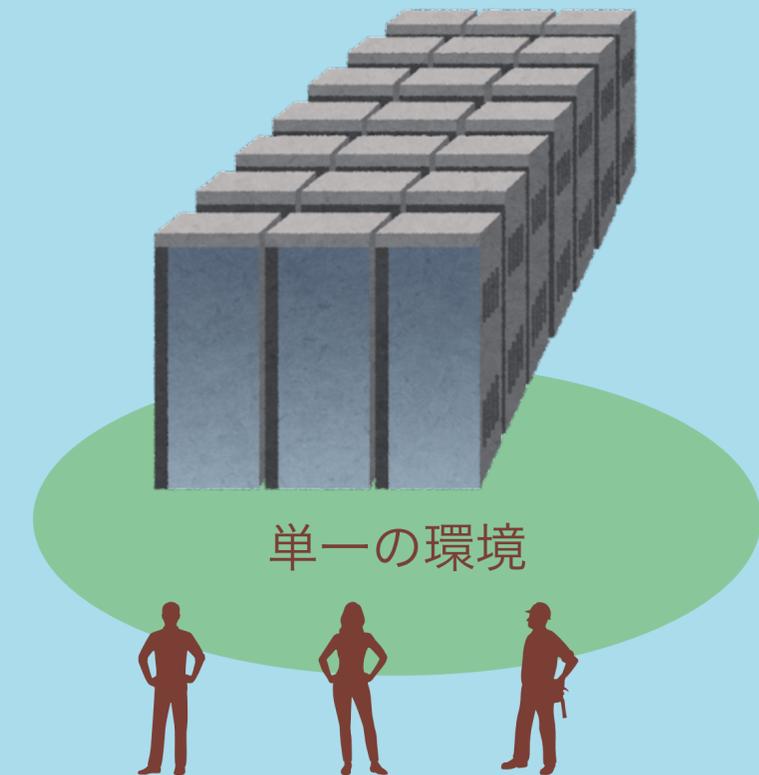
mdx ≠ スパコン

スパコン

- 大規模・高速計算向け
- 計算目的に応じた最適化
 - + 管理者による単一の環境
 - + スパコン毎に得意・不得意な分野がある

mdx

- 高性能汎用計算機
 - ユーザ毎に仮想化環境を提供
 - + 自身で OS や開発環境を構築する
- => 分野を超えたデータ活用へ



mdx の料金体系

資源種別	プロジェクトタイプ	利用資源	消費ポイント	消費開始のタイミング	
計算資源	通常プロジェクト	CPUパック *1 (1パックあたり)	起動保証仮想マシン用予約分 *8	0.2ポイント/時間	申請が承認された時点
			仮想マシン起動分*9	0.2ポイント/時間	仮想マシンを起動した時点
		GPUパック *2 (1パックあたり)	起動保証仮想マシン用予約分 *8	50ポイント/時間	申請が承認された時点
			仮想マシン起動分*9	50ポイント/時間	仮想マシンを起動した時点
	ノード占有プロジェクト	汎用 (CPU) ノード (1ノードあたり) *3	60.8ポイント/時間	申請が承認された時点	
		演算加速 (GPU) ノード (1ノードあたり) *4	800ポイント/時間		
ストレージ資源	仮想ディスクストレージ (1GBあたり) *5	0.03ポイント/日	申請が承認された時点		
	高速内部ストレージ (1GBあたり) *6	0.03ポイント/日			
	大容量ストレージ (1GBあたり) *6	0.02ポイント/日			
	オブジェクトストレージ (1GBあたり) *7	0.01ポイント/日			
グローバルIP アドレス		無料 ※但し、申請は必要			

200 TB のストレージと
50 CPU を利用する場合、
~200 万円/yr

民間 (AWS EC2) の例
48 CPU: ~15000 \$/yr
ストレージ (EBS)
16 TB (最大): ~900 \$/month

1 ポイント = 1 円

mdx ポイントが補助される研究公募: JHPCN 等

オンプレ or クラウド

木曾観測所 (on-premise) の抱える問題点

- ネットワーク帯域 (SINET で 10 Gbps 化)
- 物理的・電力 (熱) 的制約
 - + ドーム、および本館計算機室はすでに利用率高
- 管理コスト (観測所の人員)



木曾観測所 本館計算機室

mdx (cloud)

- SINET で日本全国に高速接続
- 障害対応や拡張性
 - <-> 利用料の発生 (ただし民間の 1/10 程度)



天文とビッグデータ

近年のサーベイではデータサイズが爆発的に増加

=> データアーカイブや解析に対する
データセンターの役割が重要

- PFS / Subaru
国立天文台内にデータ解析用プラットフォームを設置
- LSST (Rubin Observatory)
Interim Data Facility (IDF) on Google Cloud

Tomo-e Gozen / Kiso Schmidt

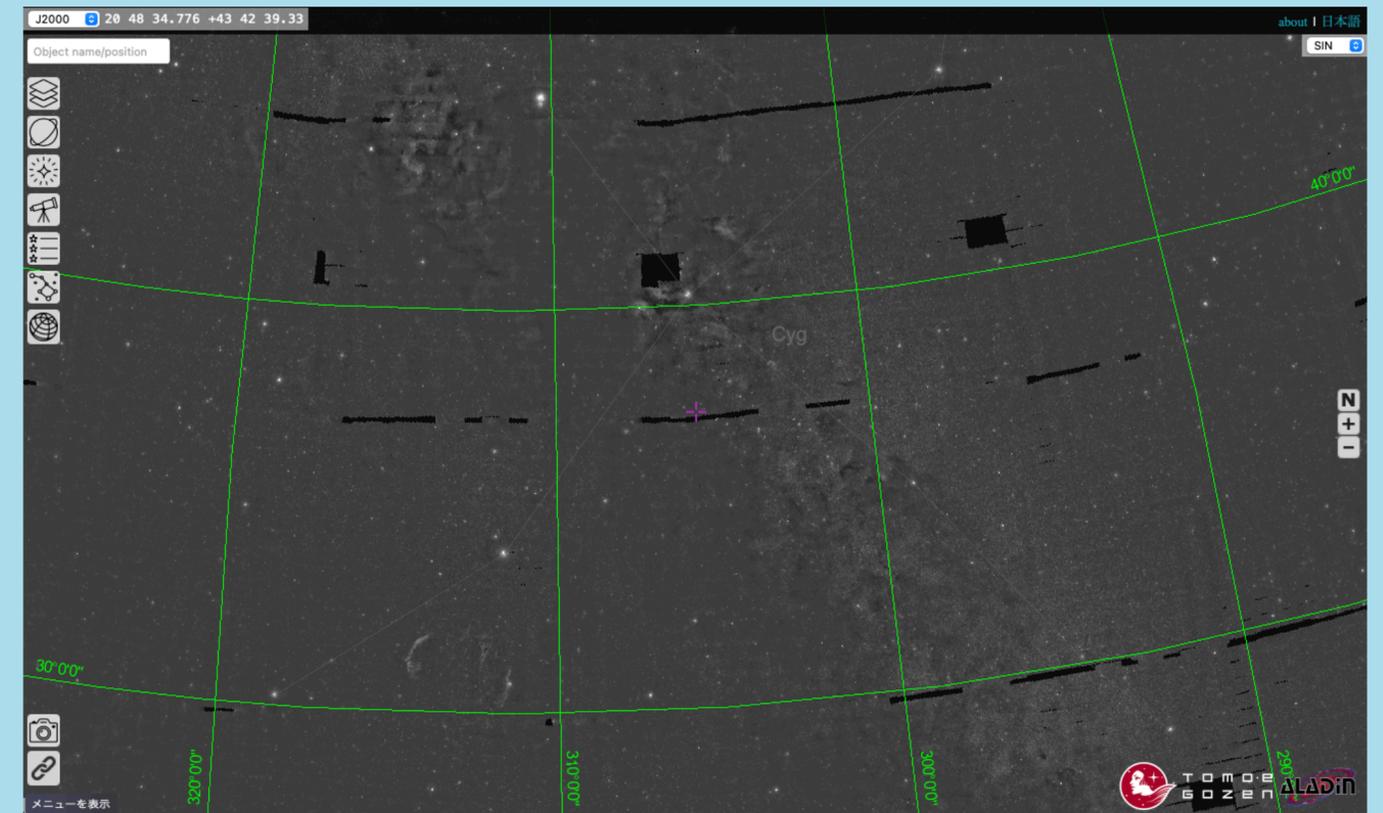
=> mdx 上への展開を目指す

Sky Surveys: Data Volumes

Sky Survey Projects	Data Volume	
DPOSS (The Palomar Digital Sky Survey)	3 TB	1990s
2MASS (The Two Micron All-Sky Survey)	10 TB	
GBT (Green Bank Telescope)	20 PB	2000s
GALEX (The Galaxy Evolution Explorer)	30 TB	
SDSS (The Sloan Digital Sky Survey)	170 TB (DR15) 40 TB	2010s
SkyMapper Southern Sky Survey	500 TB	
PanSTARRS (The Panoramic Survey Telescope and Rapid Response System)	~ 40 PB expected	ZTF: ~ 1 PB/yr
LSST (The Large Synoptic Survey Telescope)	~ 200 PB expected	2020s
SKA (The Square Kilometer Array)	~ 4.6 EB expected	(from Zhang 2015)

mdx への展開

1. mdx 上にデータアーカイブを構築
 - (ほぼ) リアルタイムでのデータ転送
 - データベース再構築
 - => obslog としての利用に最適化
 - Sky Atlas の機能拡張
2. 「秒」スケールの測光データベース
 - NEO 検出アルゴリズムを流用 (現状は捨てられているデータ)
 - 突発天体 (異常) 検知システム => データ発信へ



「分野を超えたデータ活用」が mdx には求められる