Tomo-e Gozen data platform on mdx

Satoshi TAKITA (loA/UT)

about mdx



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 (学術版の AWS のようなもの)

実体は東京大学柏川キャンパス

利用料金が発生 民間に比べると 1/10 程度 「JHPCN 共同研究課題」を通した研究費補助 「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ホーム。

Ӛ データ活用社会創成プラットフォーム 3本柱

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、データブ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(仮想データブラットフォーム:仮 想DP)を、使いたいときに即時に構築する

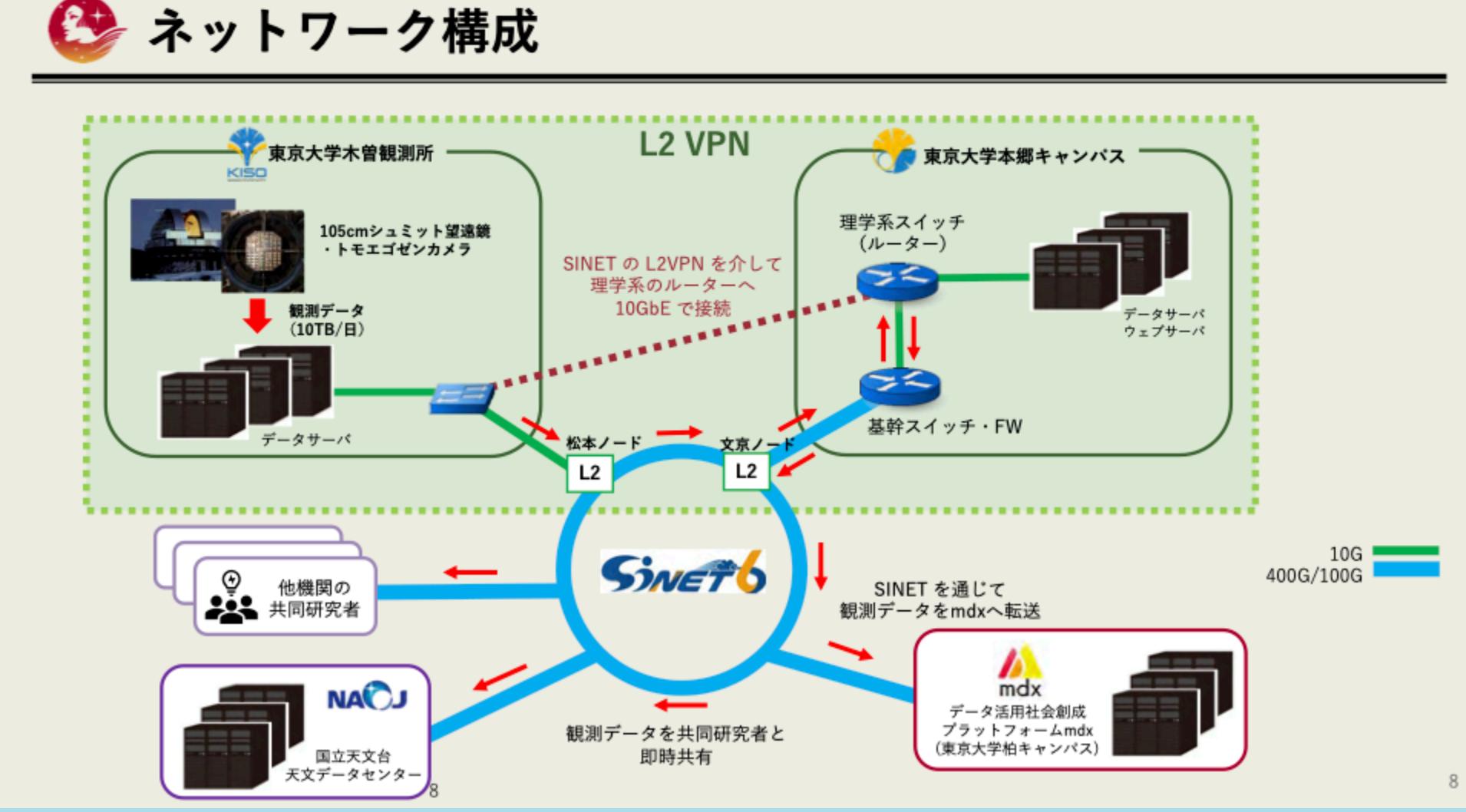
高性能計算環境によるデータ科学と計算科学の融合 データ科学、計算科学の手法を融合し、さらに国内最高の計算環境を用いて他に無い 高精度の予測を行えるようにする

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 ティーを形成し、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 and kiso-SINET



project mdx

背景

- 木曽観測所の計算機群の問題
- +物理的、電力(熱)的制約
 - "実効的に" 利用可能なストレージは 1 PB 程度
- + 管理コスト (観測所の人員)

目的

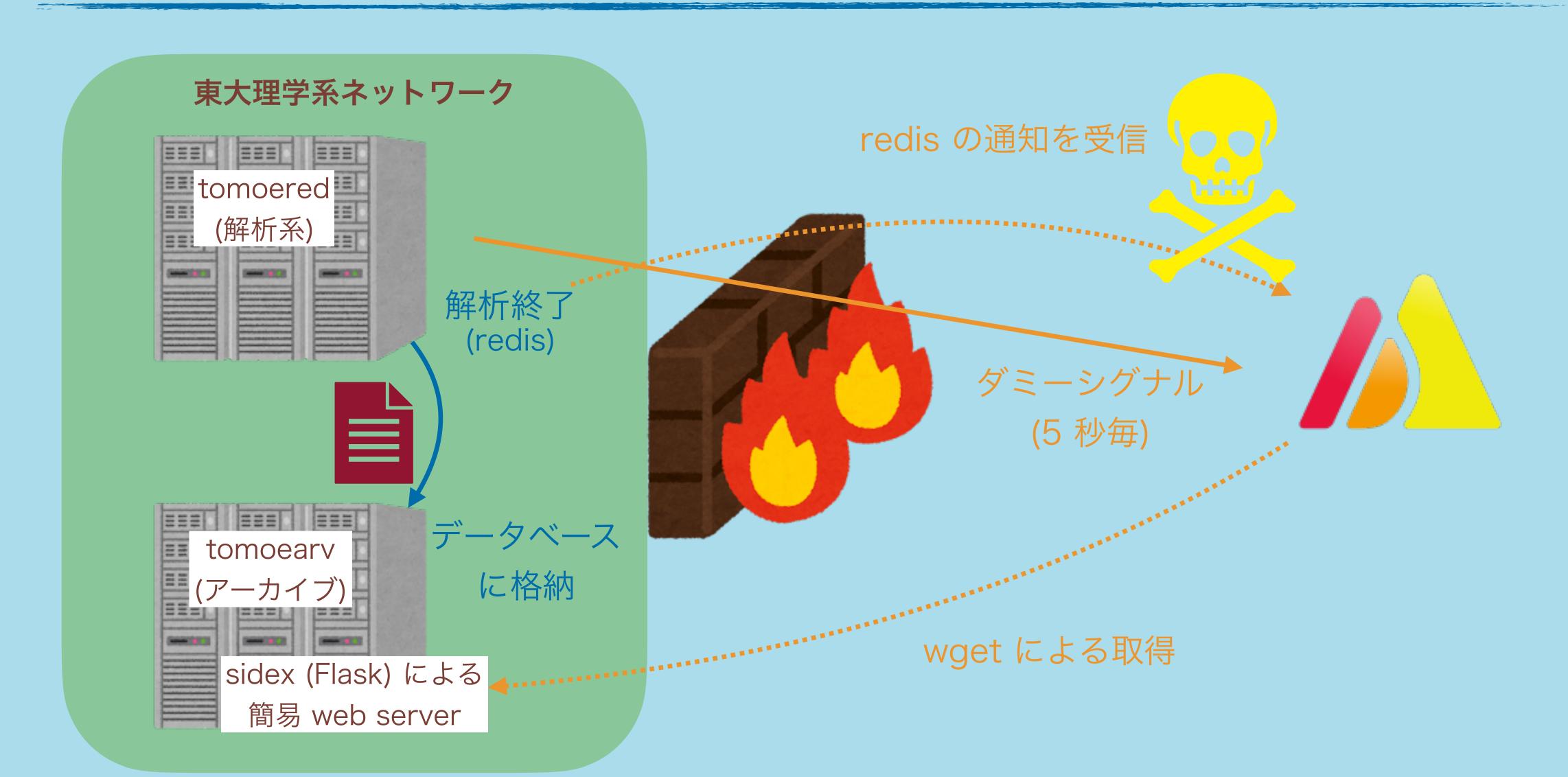
- データ解析機能の一部を木曽観測所から mdx ヘオフロード
 - 1. データアーカイブ
 - 2. ライトカーブデータベース
- mdx を起点としたデータ発信

1. alternative data archive

データアーカイブの作成

- スタック済み画像のリアルタイム転送
 - + Tomo-e のデータ解析で利用されている redis 通知
 - + sidex (Flask ベースの簡易 http server) 越しの wget
- データベース
 - + Apache Arrow 形式 (列指向データベース)
 - + pyarrow/duckdb
- web UI
- + 現状のアーカイブを流用・強化

data transfer



database

観測ログ

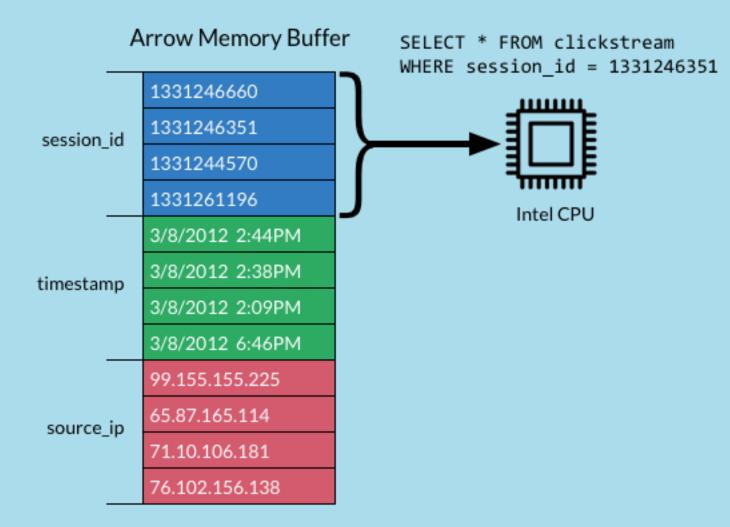
- 観測システムとは独立
- 毎晩の観測終了後に一度だけ追記
- 検索性能が重要
- => OLAP 的なデータベースが有効

Apache Arrow

- 列指向データベース
- parquet によるファイルベース
- pyarrow や duckdb 等のライブラリ
- (PG-Strom で PostgreSQL に組み込み可)

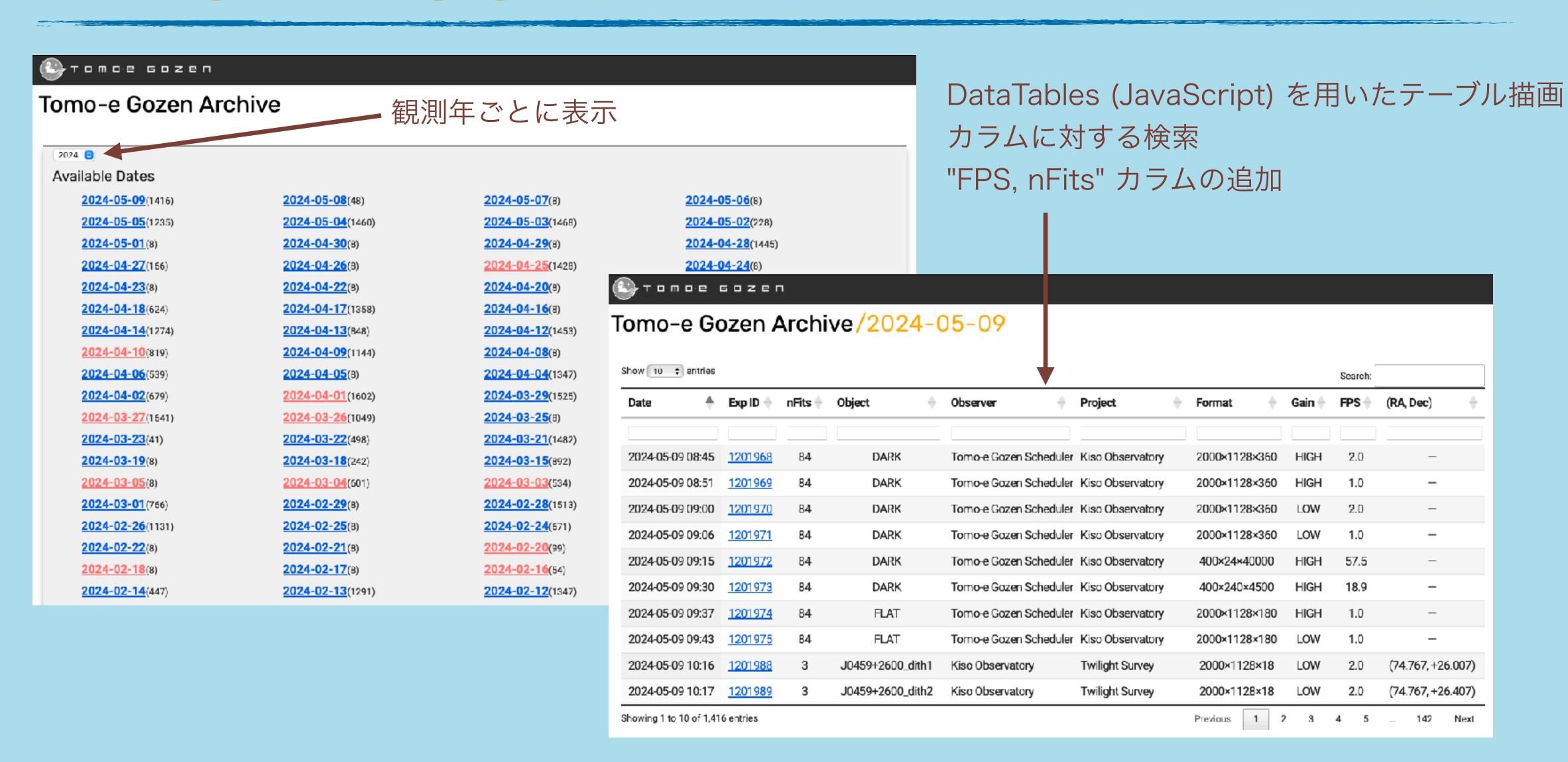
		session_id	timestamp	source_ip
Row	1	1331246660	3/8/2012 2:44PM	99.155.155.225
Row	2	1331246351	3/8/2012 2:38PM	65.87.165.114
Row	3	1331244570	3/8/2012 2:09PM	71.10.106.181
Row	4	1331261196	3/8/2012 6:46PM	76.102.156.138





https://arrow.apache.org

interface

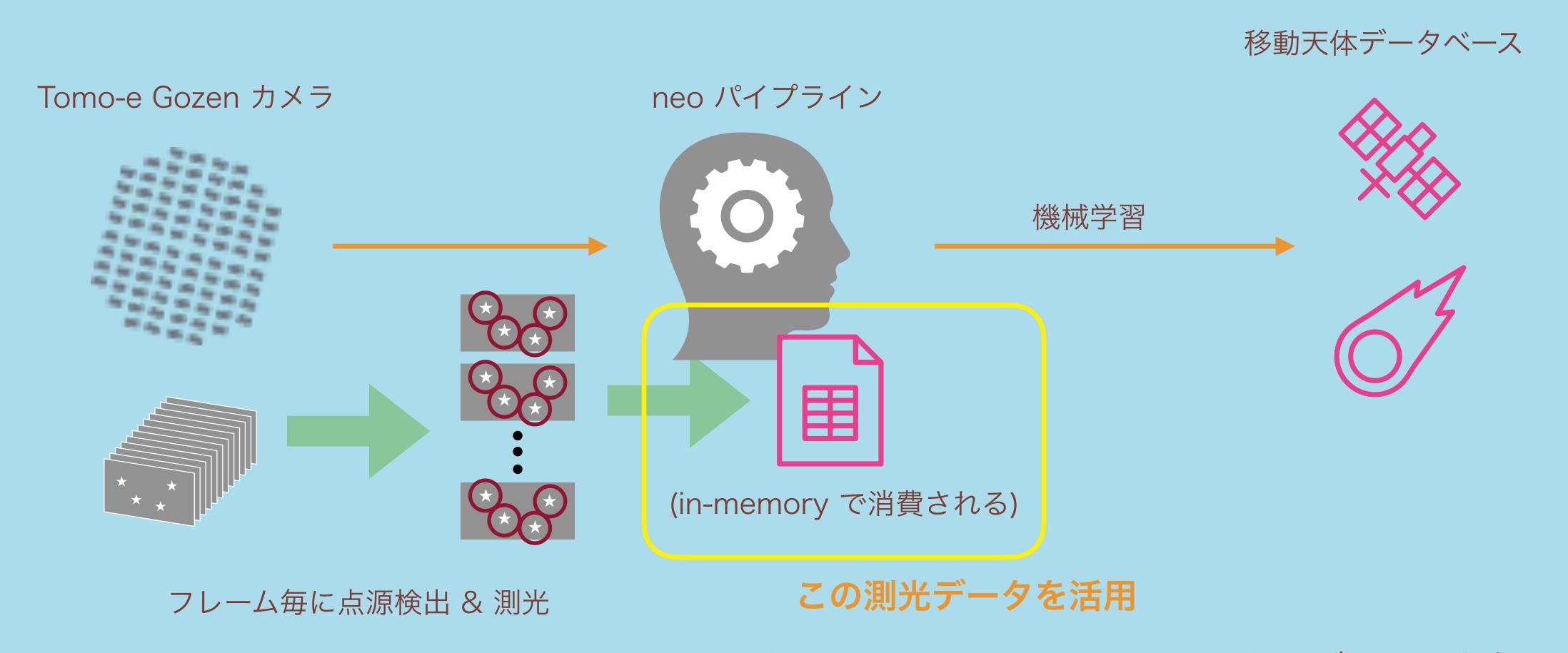


2. light-curve database

Tomo-e が誇る「秒」スケールのライトカーブ

- neo / 高速移動天体検出パイプラインによるフレーム単位での測光データ
 - +6 or 9 秒のライトカーブ
- Gaia データとのマッチング
- + 数時間や日を跨いだデータの結合
- + flux calibration
- データ発信
 - + skyatlas との連携
- + デブリデータ

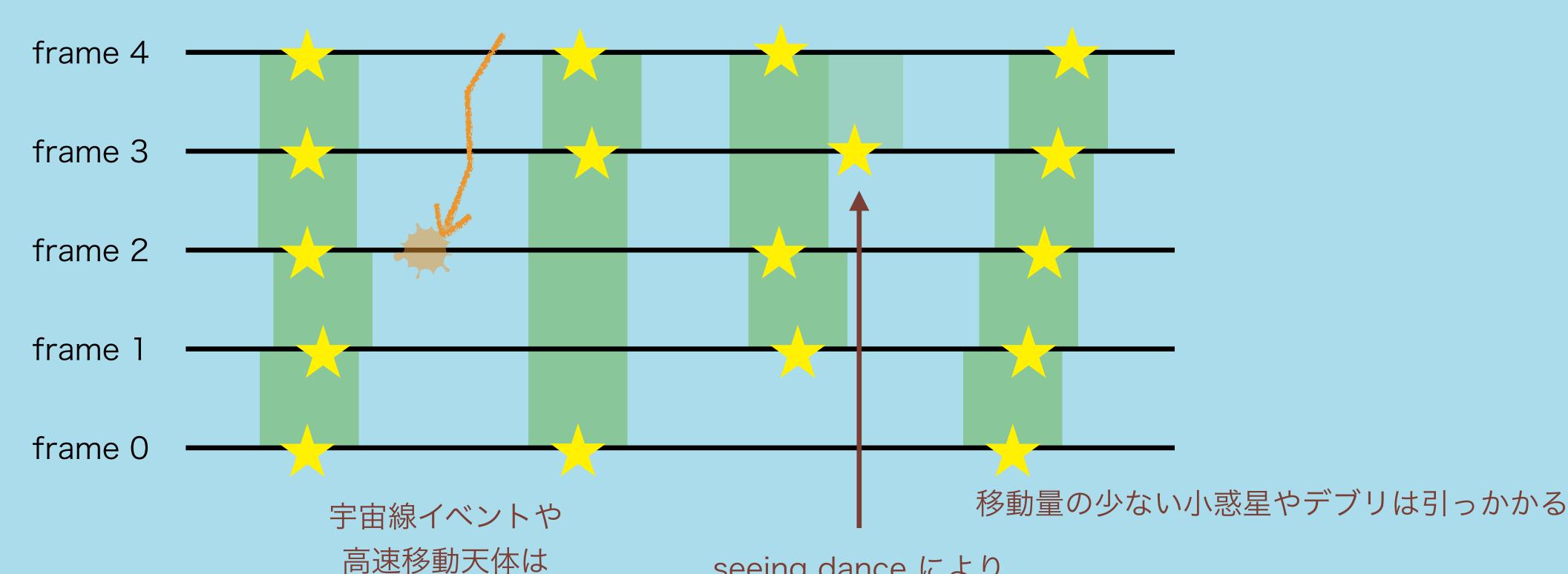
neo data



試験として 1 exposure (84 fits) 分のデータを出力

data grouping / matching

画像上の x, y 座標 -> +- 3 px でマッチング 直近の検出位置を基準 => Gaia DR3 と 1 arcsec でマッチング

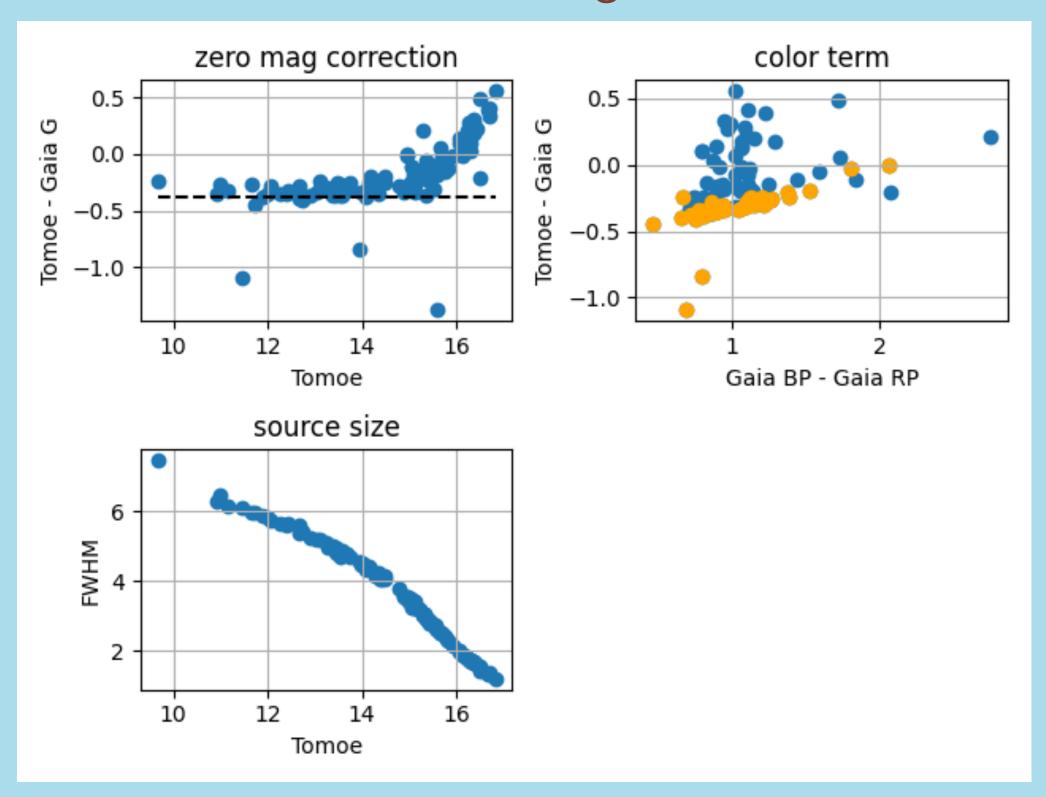


単発のイベント

seeing dance により 外れ値がでることもある

flux calibration

Gaia との比較により zeromag, colour-term を導出



現状の問題点

- Tomo-e データは
 sep (SourceExtractor) の
 magauto の値を出力
- -> 固定 aperture の測光値が必要

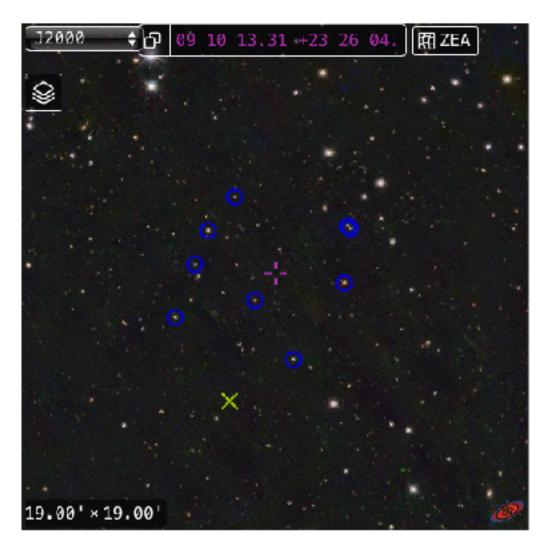


- 21 px diameter の測光値
- Gaia データから Tomo-e mag への変換



interface

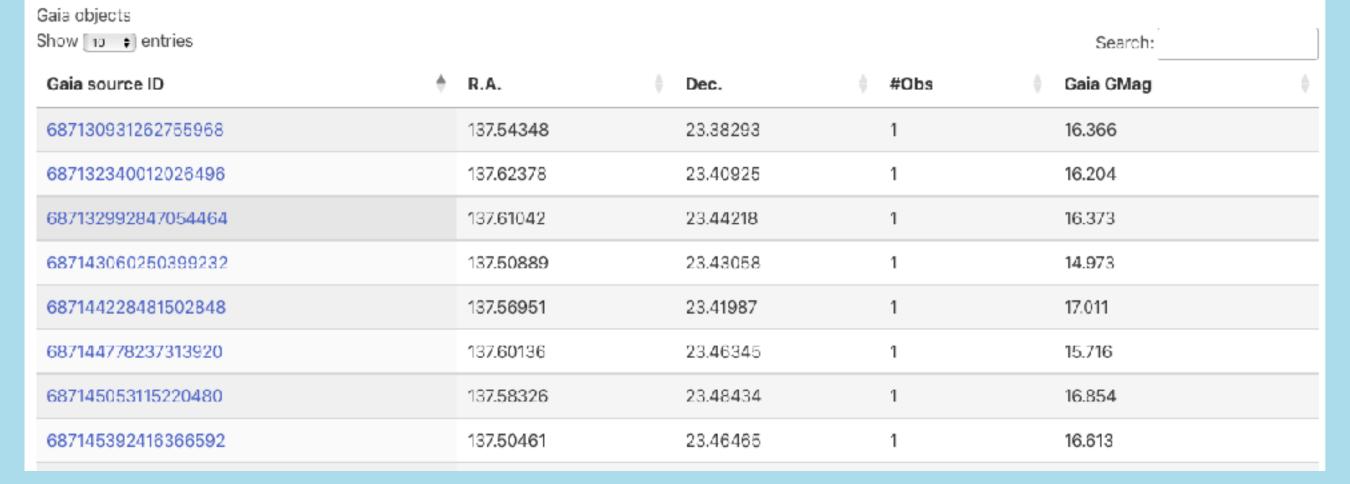
Tomo-e Gozen (Sub-)Second-scale Lightcurve

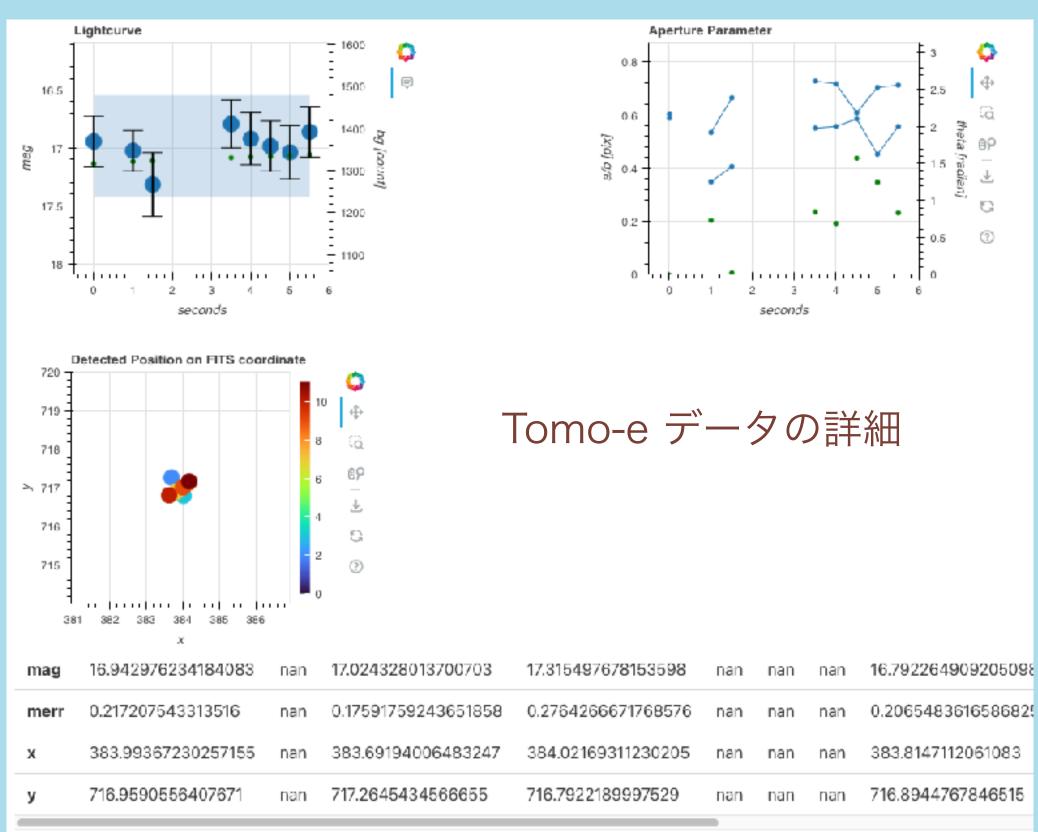




青〇: Gaia とマッチングした天体

緑×: Gaia にない天体





issues to do

データアーカイブ

- ブロックリストの対応
- ストレージ容量問題

利用料金の問題で mdx のみでは困難

=> HPCI 共有ストレージ利用研究課題

ライトカーブ

- データの準備
- + neo パイプラインの改修・常駐化
- データベースの検索が遅い
 - + Tomo-e と Gaia の JOIN が問題
 - + 恐らく parquet (ファイルベース) で扱っているのが原因
 - => Tomo-e 側に Gaia の情報 (一部) を加えるべきか

summary

mdx 上に新たなデータプラットフォームを開発中

- JHPCN 共同研究課題や HPCI 共有ストレージの利用
- SINET による 100 Gbps ネットワーク

データアーカイブの機能強化・「秒」スケールライトカーブ

- skyatlas との連携
- 動画データの有効活用
 - + 産学官連携
 - => 分野を超えたデータ活用へ