



Searching Large Scientific Data

John Wu



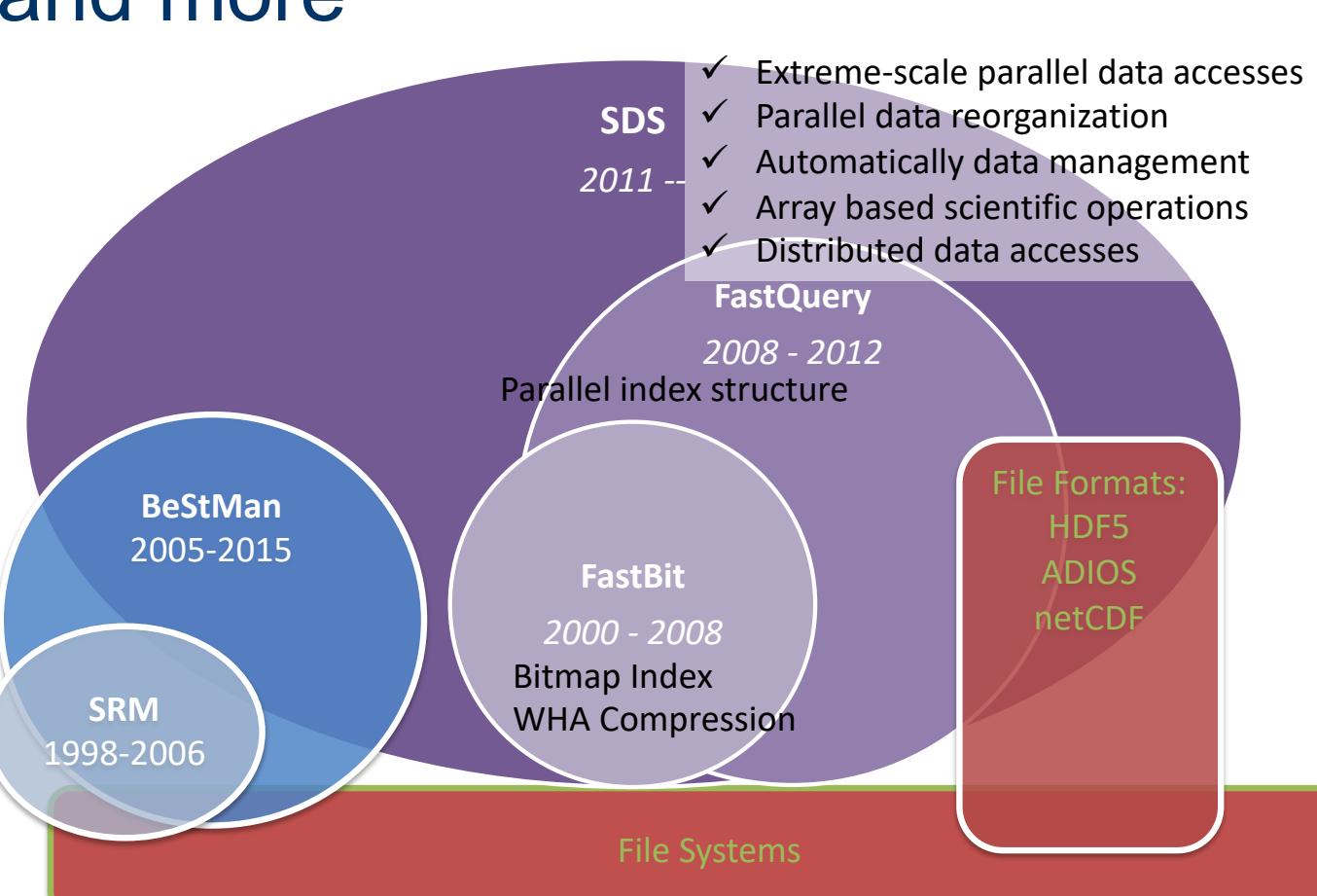
Scientific Data Management Group



Overview

Data management is key in simulation and analysis for speed, memory and storage efficiency

- Capability:** Premier group in algorithms and data structures, for data management, I/O and storage
- Impact:** Improve applications performance, enable new science problems, more productive science
- Stakeholders:** Application scientists and facilities; all SC, EERE, Health, and more



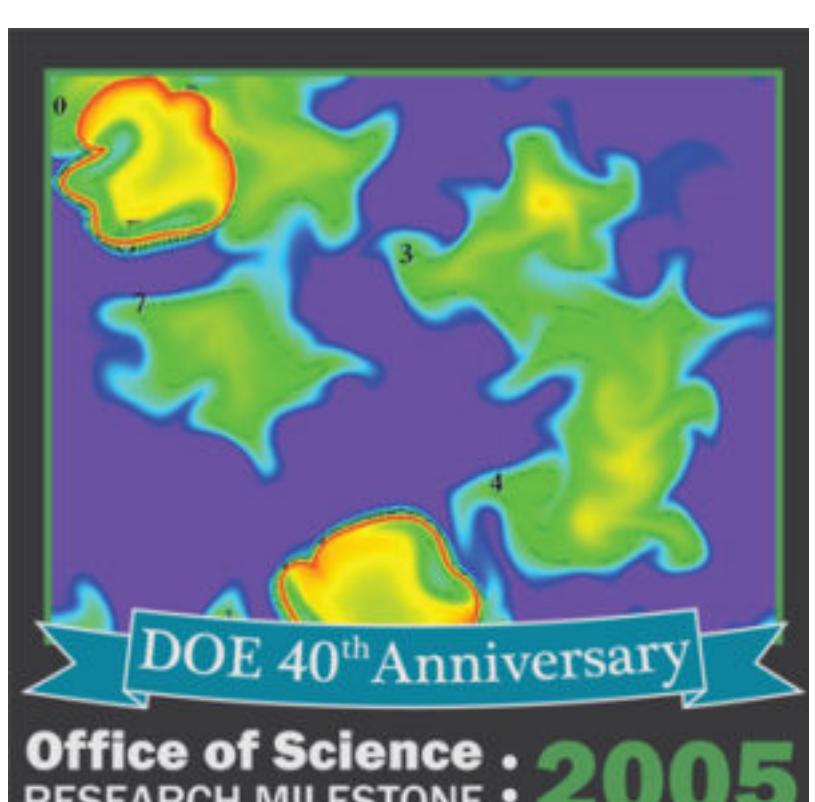
FastBit Indexing

FastBit is extremely efficient in many applications: high-energy physics, combustion, astrophysics, network security, drug discovery, ...

The efficiency comes from new methods and algorithms, careful software engineering, and rigorous theoretical analyses to prove optimality

Won R&D 100 Award in 2008

Selected to be a part of the DOE40 celebration as a one of ASCR research milestones



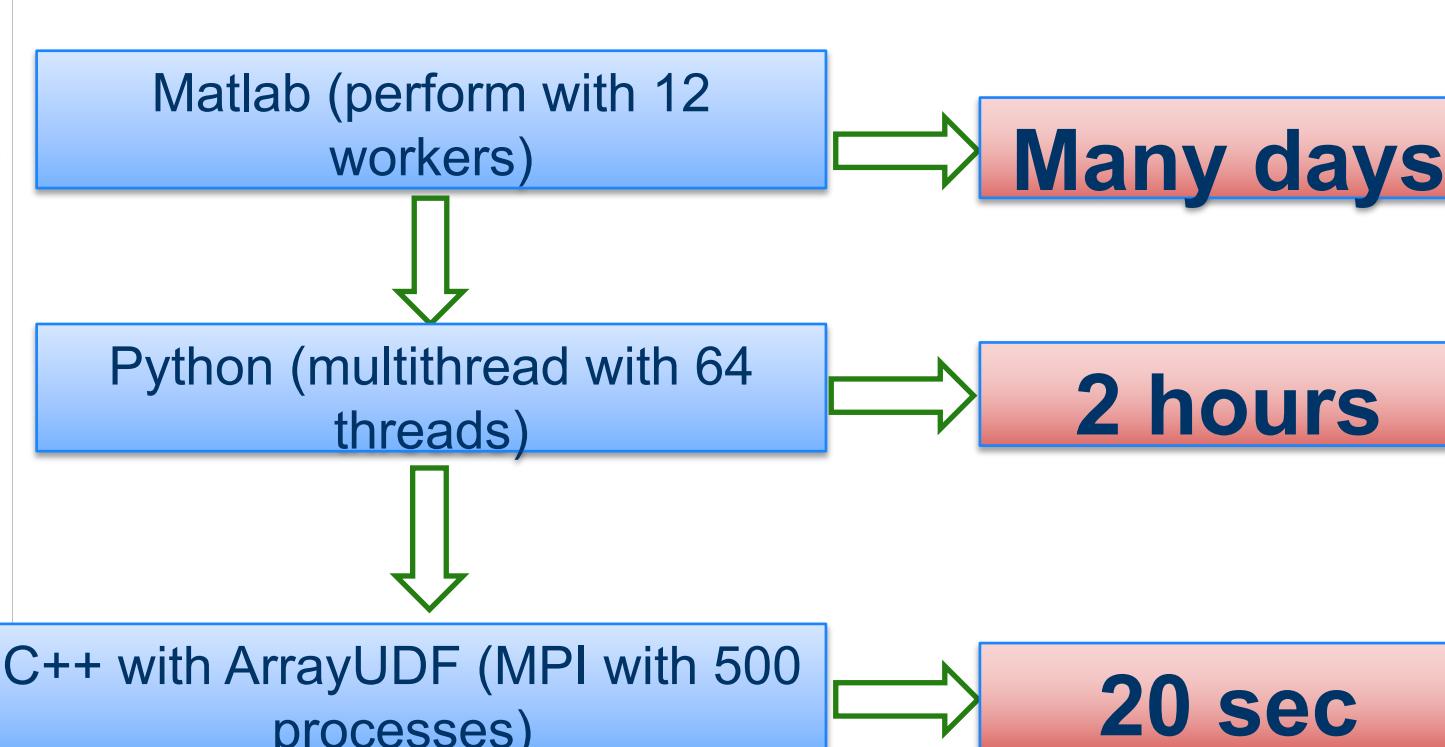
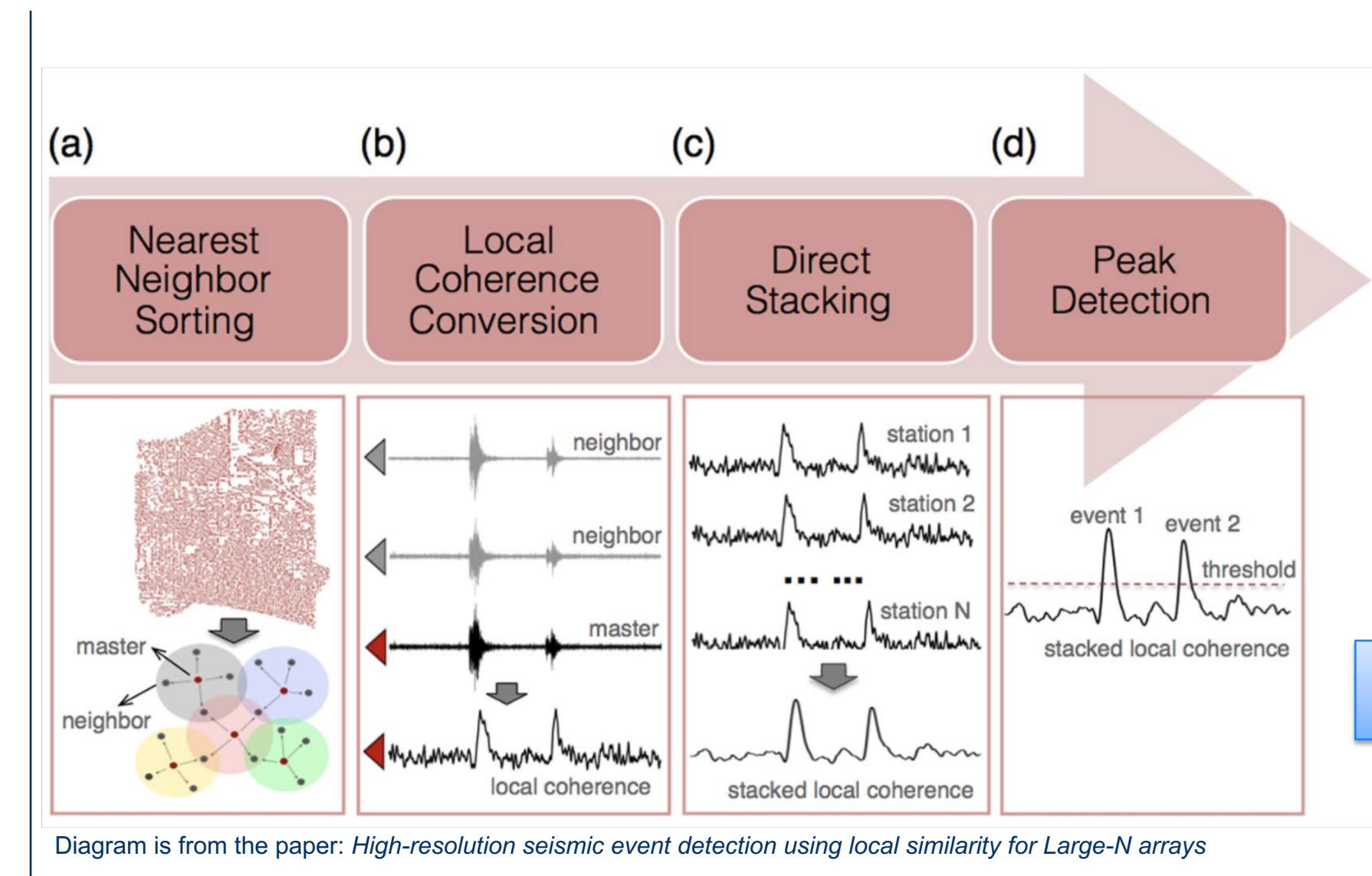
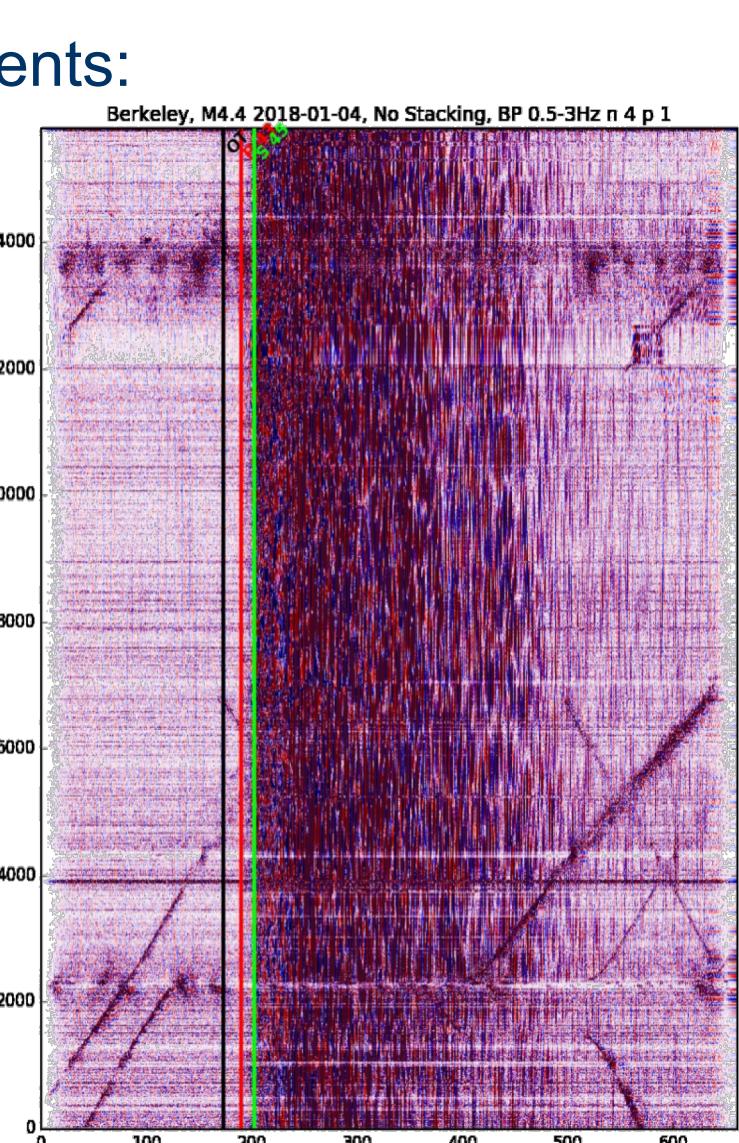
Applying Scientific Data Services to Distributed Acoustic Sensing Data

Distributed Acoustic Sensing (DAS) collects seismic data at high resolution, and requires high-performance analysis capability!

Detect different seismic events:

- earthquake
- vehicles
- ...

Much more sensors
Big data
Algorithm Needed
Better event detection



Array Data Model

Example task: Join data from different sources is useful but time-consuming, e.g., comparing new and old astronomical observations to identify supernovae. Completing this step quickly is critical before the transients disappear.

Our Solution:

Define a theoretical framework to identify optimal approaches [SIGMOD 2016]

Create incremental strategy to process periodic additions [SIGMOD 2017]

Design a distributed cache for query processing [SSDBM 2018]

Application:

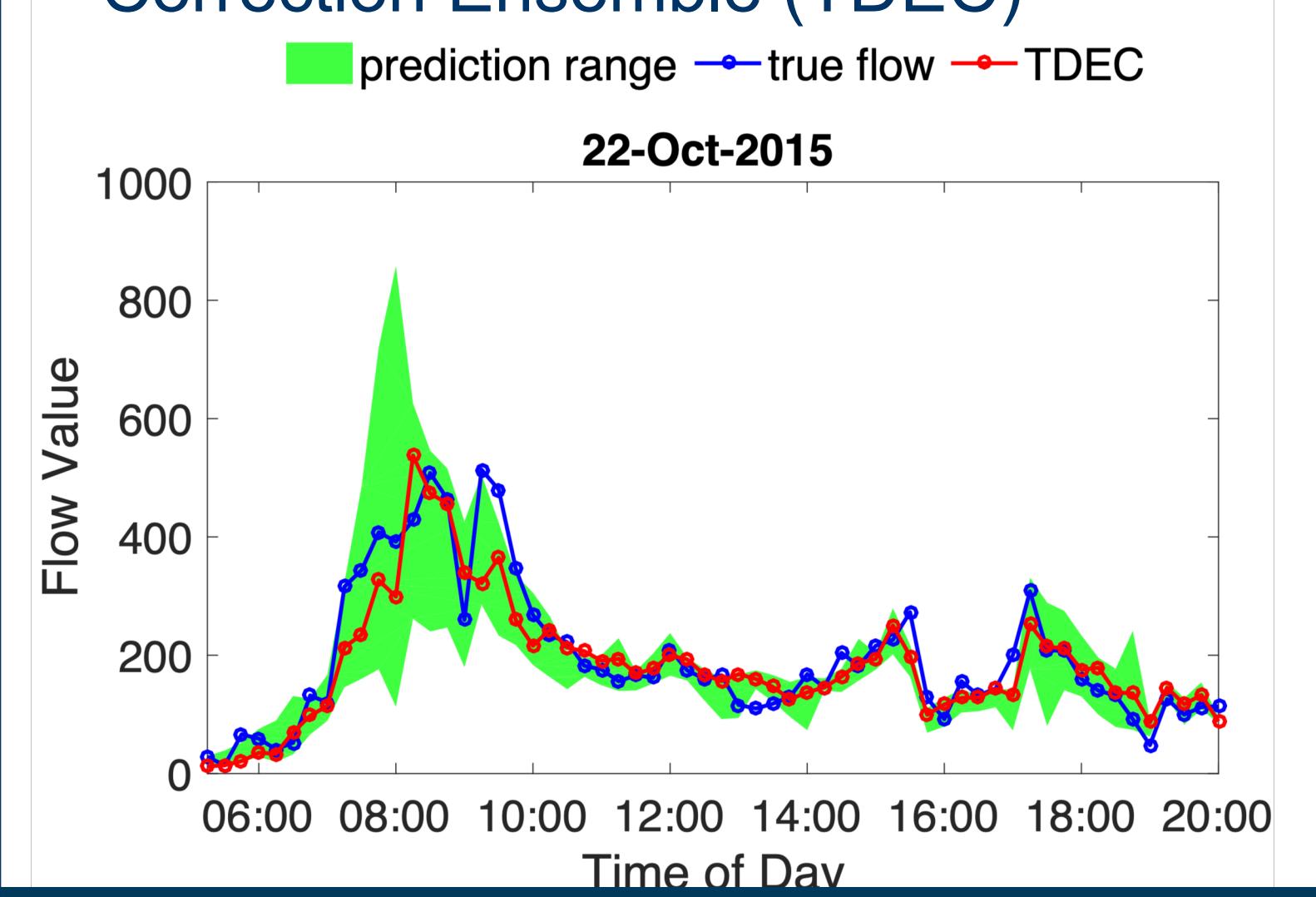
Help locate the past observations of two neutron stars in a collision (joint authorship in a Science paper).

Streaming Data

Example Task: one-hour-ahead arterial traffic flow (#cars/hour) prediction in Arcadia, CA, at 15 min granularity (15 min/step)

Base models: (1) Autoregressive moving-average, (2) Partial Least Squares, (3) Support Vector Machine, (4) Kernel Ridge Regression, (5) Gaussian Process Regression

Our method: Time Decay Error-Correction Ensemble (TDEC)



Additional Use Cases

Our technology significantly accelerates scientific applications

Application	Approach	Size	Speed
VPIC	Block Index	100x	5x
Plasmas	Block Index	100x	5x
Cosmology	AMR-Index		500x
AMR	AMR-Index		
Brain EEG	Statistical Similarity	106x	
Power Grid	Statistical Similarity	198x	
Mass Spec	Multilayout		90x

