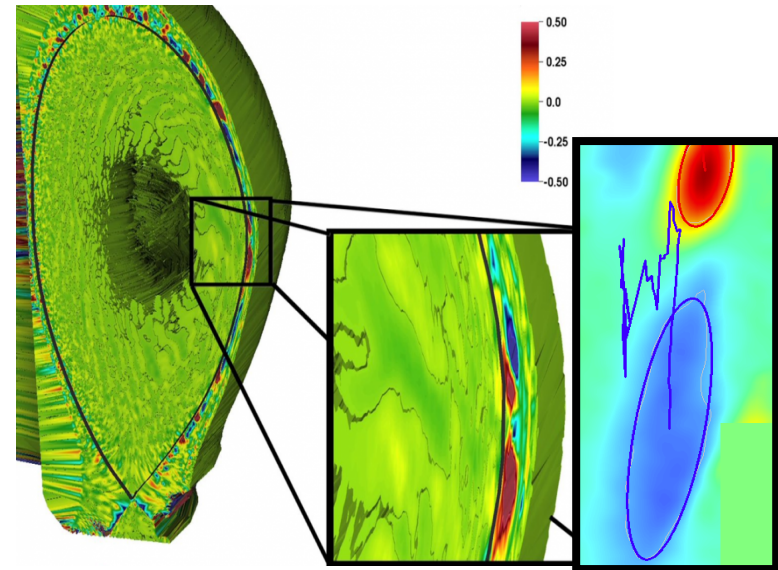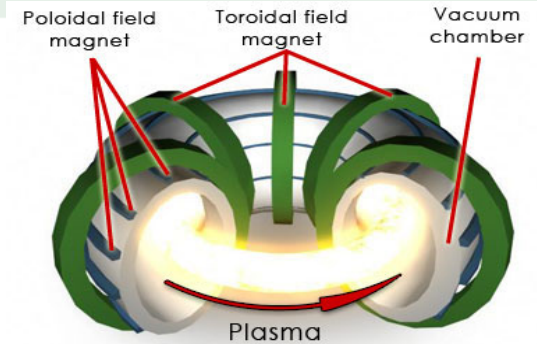# Technology for Distributed Streaming Analytics

**John Wu**
**LBNL**

# Use Case 1: Near Real-Time Feature Detection

➢ Fusion experiments are conducted at centralized facilities

  ➢ Junior researchers often operate the devices, while senior researchers offer advices from afar

  ➢ There are 10s of minutes between runs/shots

➢ Need for distributed analysis

  ➢ The experimental facility may not have enough computing power

  ➢ Need to compare experimental measurements against simulation predictions

  ➢ Measurement data ~GB/s, simulation data ~TB/s, need significant computing power for analysis

➢ Distributed in transit processing

  ➢ Make more processing power available

  ➢ Allow more scientists to participate in the data analysis operations and monitor the experiment remotely

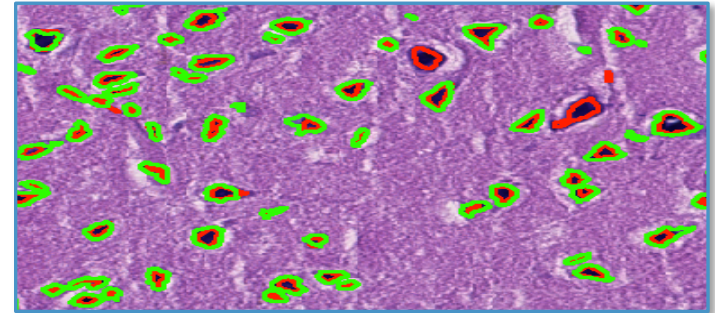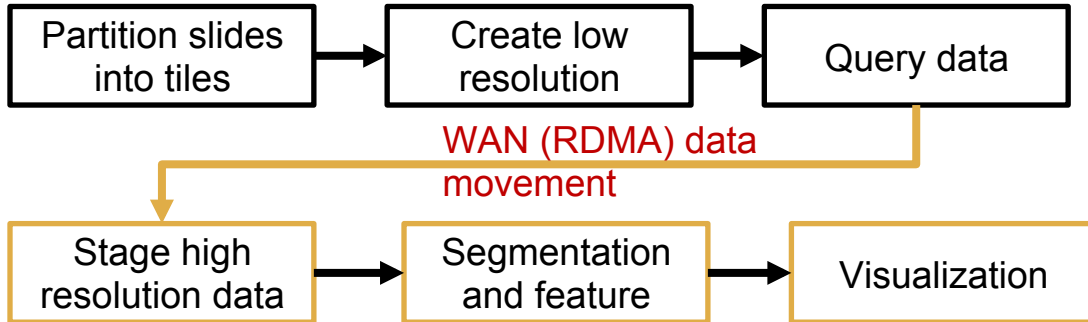  ➢ Enable scientists to share knowledge and processes

Blobs in fusion reaction
(Source: EPSI project)

Blob trajectory

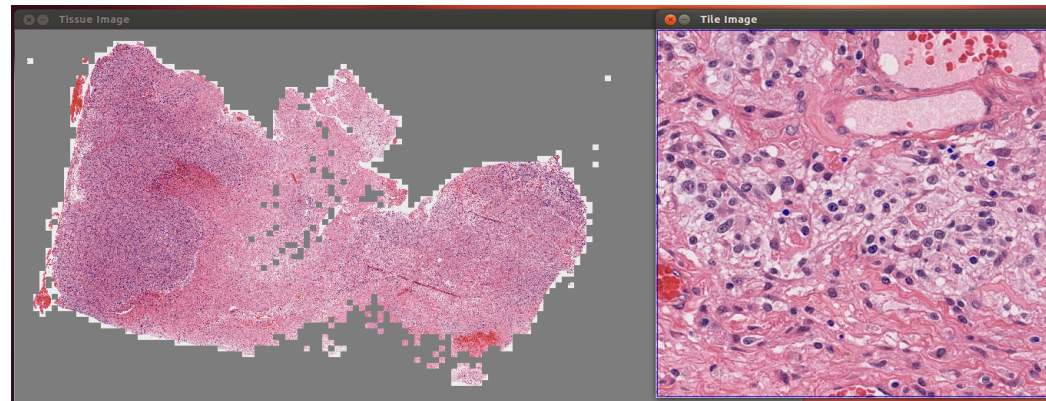Wu, Sim, Choi, Churchill, Wu, Klasky, Chang, 2014

**Challenge:** identify cancerous cells in tissue image (120Kx120K) while the patient waits
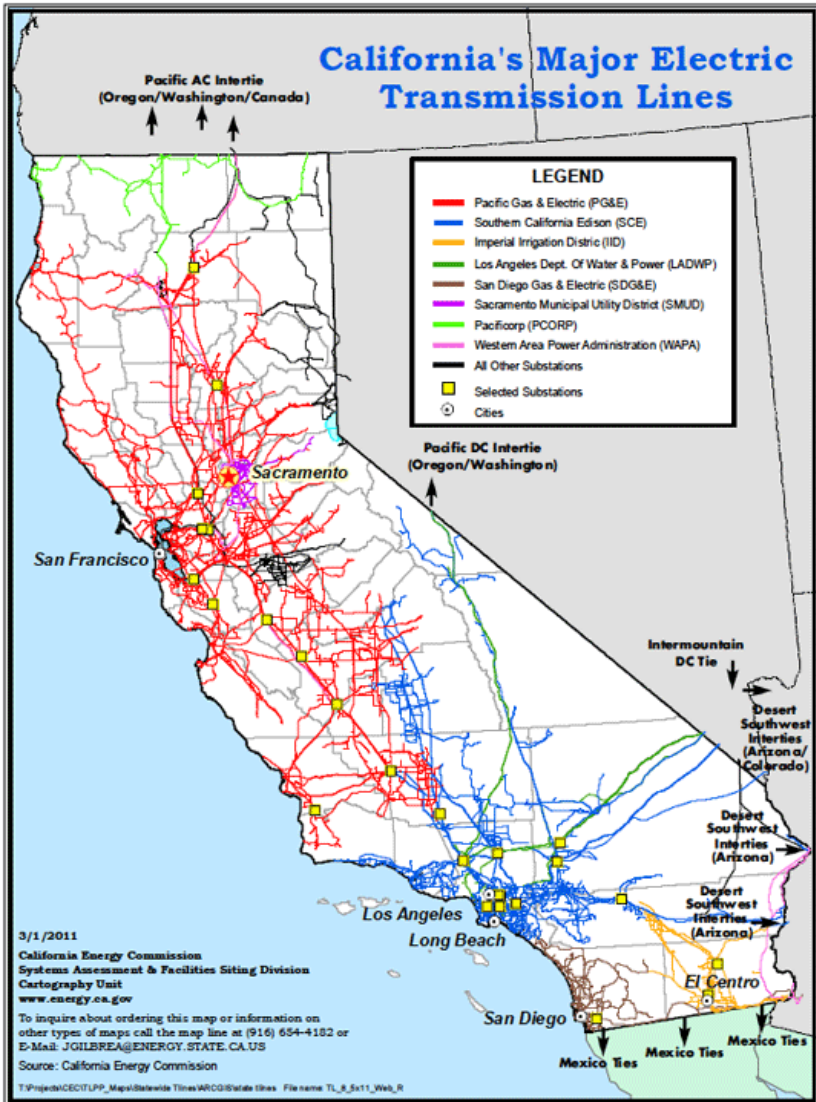


**Technologies:** (1) ICEE transport layer for wide-area, efficient transfers; (2) Longbow for very fast, low-latency connection; (3) pipelined processing on clusters

**Demo:** Tissue slides on machine in Singapore. Analysis done on cluster at Georgia Tech. Segmentation results displayed on client machine.
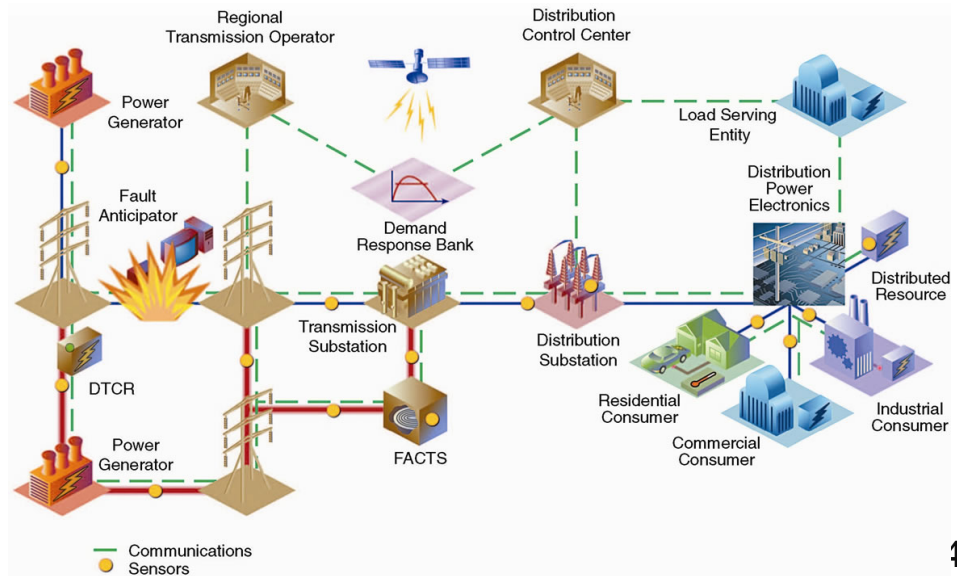
- Snapshot of adaptive processing of a remote slide
- Image broken into pieces for parallel processing
- Need to stitch the boundaries together



3

# Use Case 3: Integrate Distributed Sensor Data from Power Grid



California's Major Electric Transmission Lines

- Sensors such as Phaser Measurement Units (PMU), Smart meters, thermostats, appliances create many data streams
- Linked to other time and location-specific information (temperature, census,…)
- Proper analysis of such data is key to the vision of Smart Grid and Smart Cities

# Technology Needed for Streaming Analytics

## Velocity

- Reduce data access latency, reduce volume transferred, move analysis

## Volume

- Reduce the volume transferred, move analysis

## Variety

- Enable multiple streams of data to be analyzed together
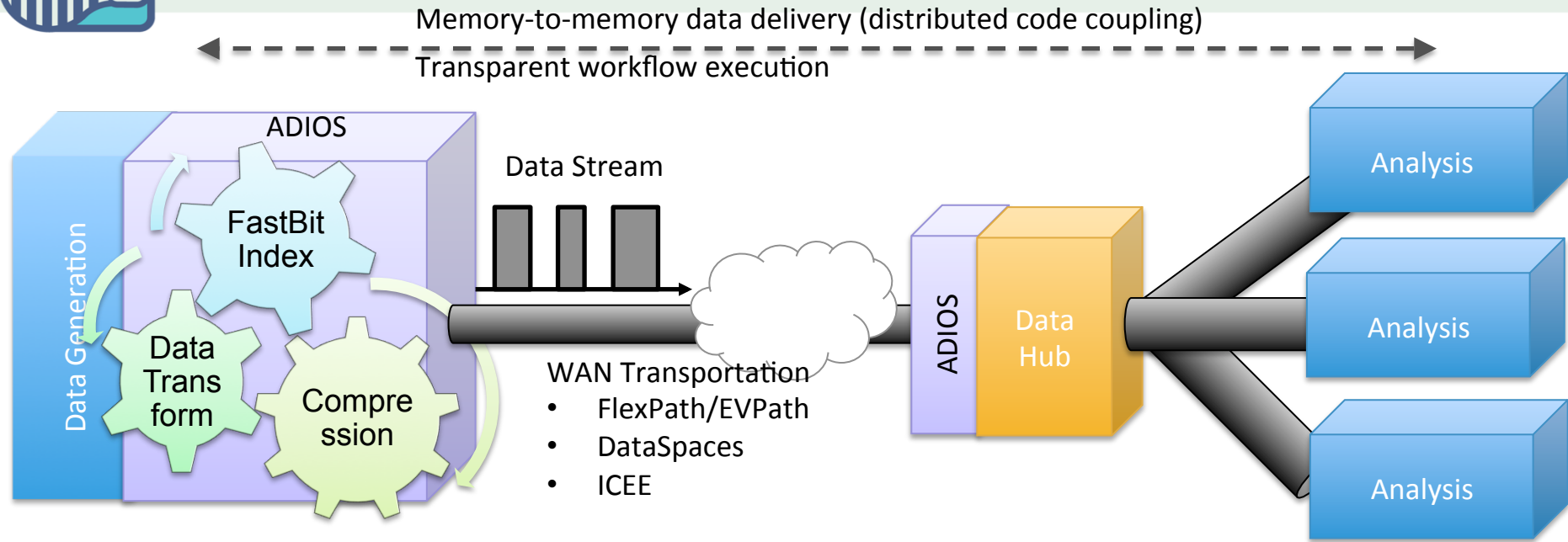
## Veracity

- Understand the trade-offs for accuracy (of the query) vs. accuracy of the results vs. performance (time to solution)

## Value

- Provide the freedom for scientists to access and analyze their data interactively

Memory-to-memory data delivery (distributed code coupling)

Transparent workflow execution

ADIOS

Data Generation

FastBit Index

Data Trans form

Compre ssion

Data Stream

WAN Transportation
- FlexPath/EVPath
- DataSpaces
- ICEE

ADIOS

Data Hub

Analysis

Analysis

Analysis

# Utilizing ADIOS in situ processing capability to keep as much of the distributed workflow in memory as possible

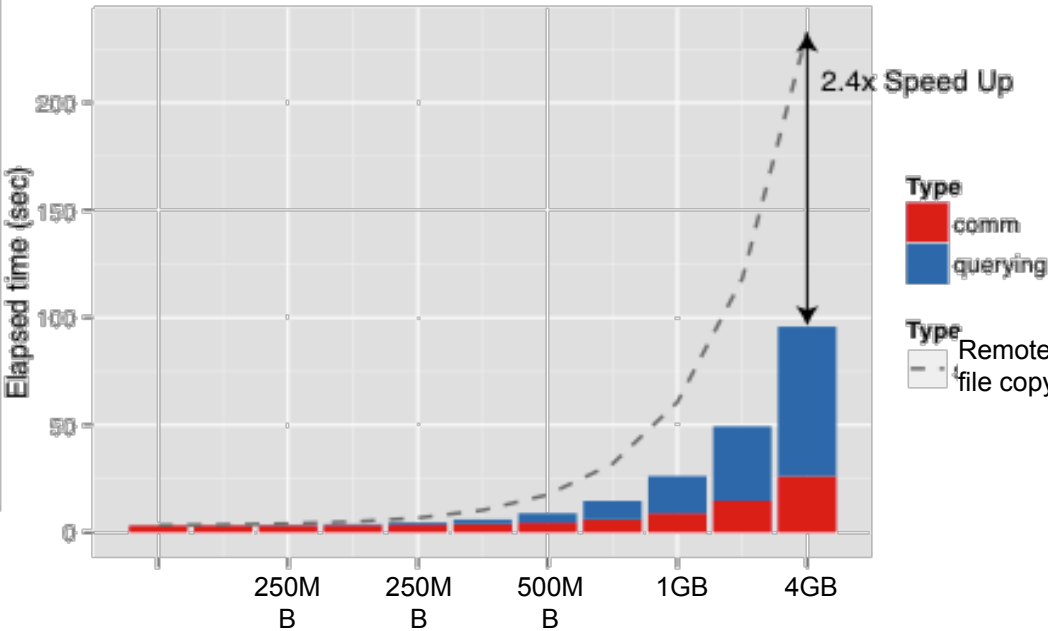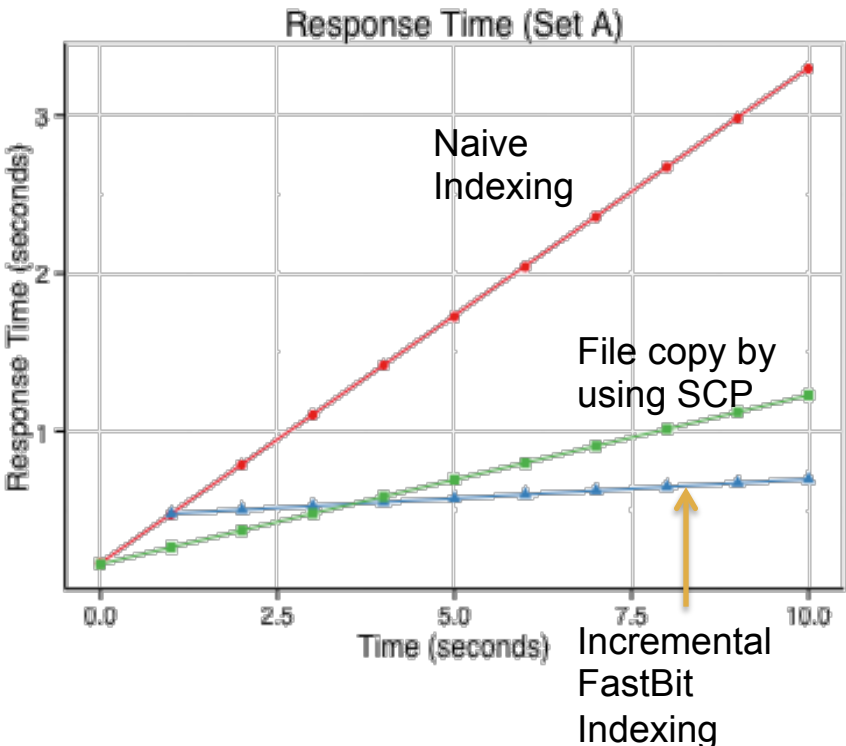- WAN transportation: FlexPath (GATech), DataSpaces (Rutgers), ICEE (ORNL/LBNL)
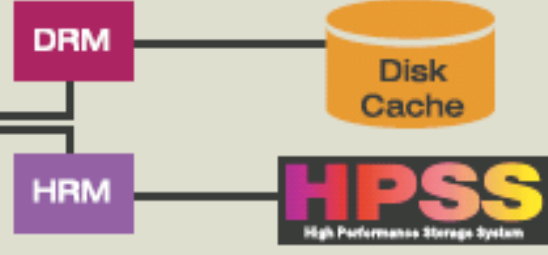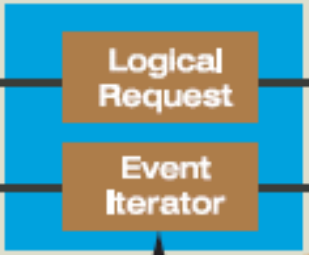
# Remote file copy VS. index-and-query

- Measured between LBL and ORNL
- Using indexes to locate necessary data, i.e., querying, reduces overall execution time
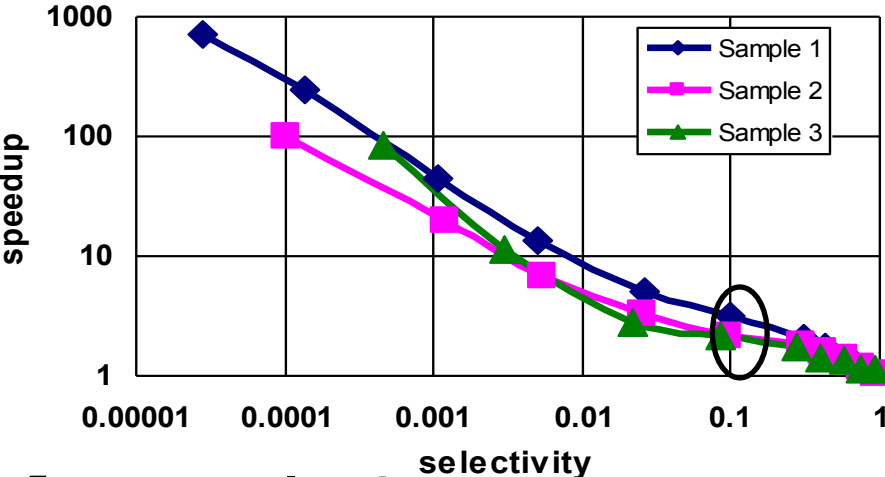


7
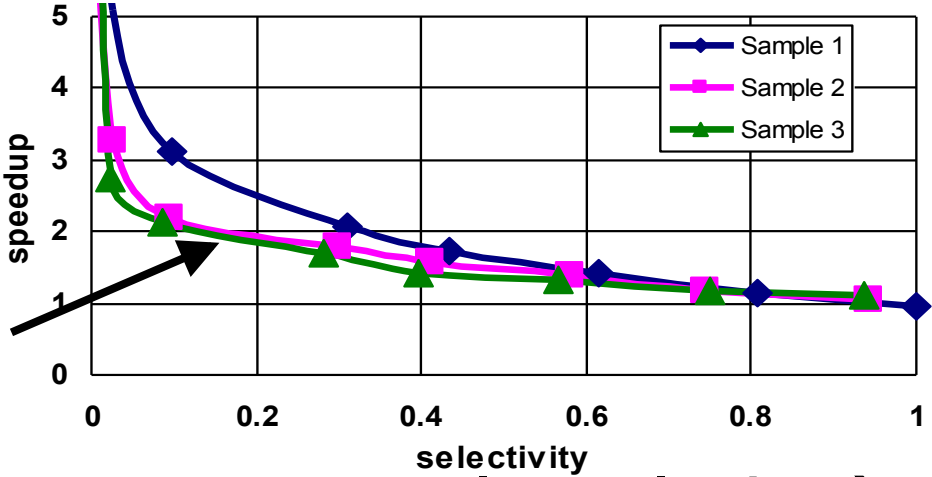
# Technology Example 4: Grid Collector



Analysis Framework — Grid Collector Servers — Remote Storage Systems
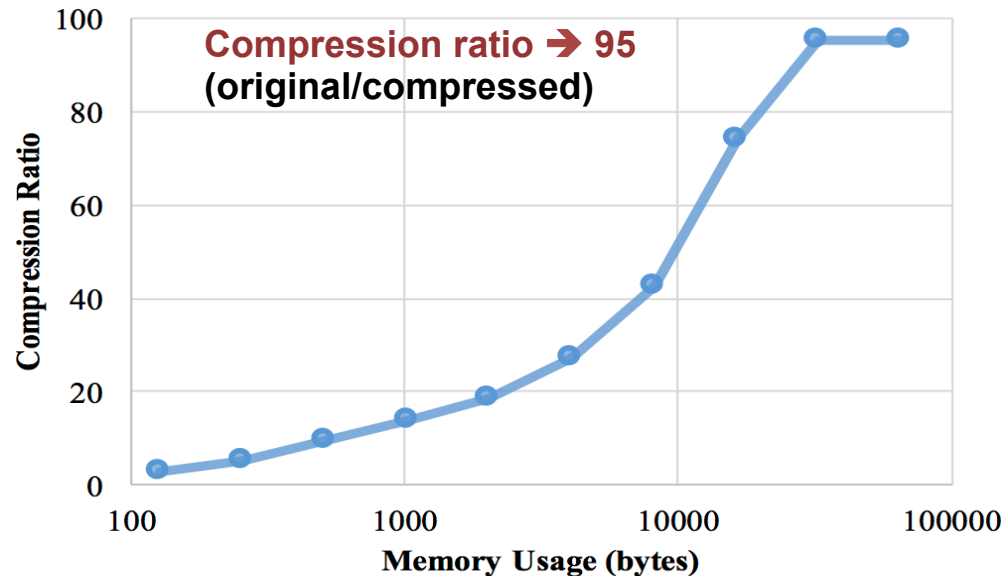


← **more selective**

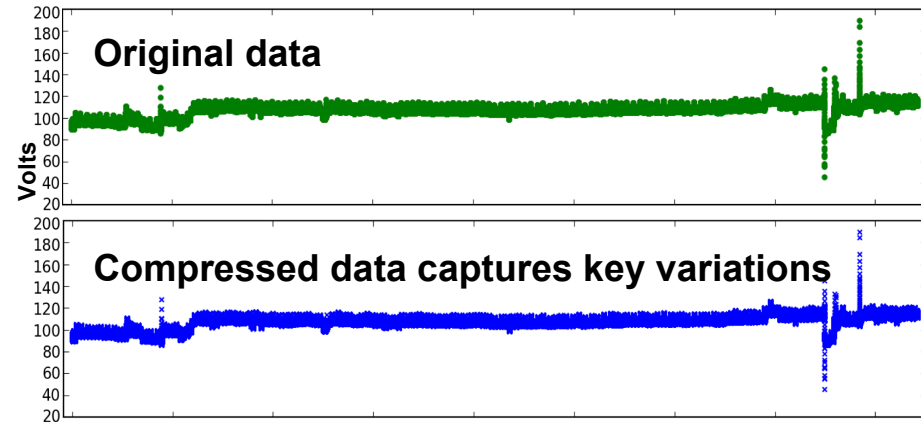**less selective** →

- Conventional compressions are based on values, but the new technique is based on Probability Density Function

- Theoretically, Locally Exchangeable Measures

- The method supports feature detection directly on the compressed data

- Test data: Micro PMU data from LBNL

- Measured data compression ratio (original size in bytes / compressed size) reaches 95, using 64KB buffer

- Compared to gzip, LEM compressed data size is under 2% of gzip-compressed data size in bytes

- Locally Exchangeable Measures, U.S. Patent pending (serial no. 14/555,365)



**Original data**

**Compressed data captures key variations**

**Compression ratio ➜ 95 (original/compressed)**

Contact: Alex Sim, SDM, CRD, LBNL <ASim@LBL.Gov>

## Algorithms

- Did not touch on algorithms for analysis, workflow orchestration, data integration, …

## Systems

- Are existing systems sufficient?

- What can be accomplished with the existing streaming systems?

## Networking needs

- Moving queries to the networking system

- QOS: guarantee delivery (because data might not be saved anywhere), guarantee bandwidth