# Division of Computational Informatics Database Group
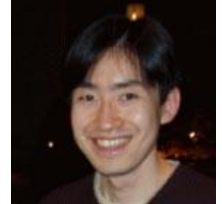
**Hiroyuki Kitagawa**
**Center for Computational Sciences**
**Graduate School of Systems and Information Engineering**
**University of Tsukuba**

# Members

## ■ Faculty

- ● Hiroyuki Kitagawa (Professor)
- ● Toshiyuki Amagasa (Associate Professor)
- ● Hideyuki Kawashima (Lecturer; Currently in HPCS Div.)
  - ■ Yasuhiro Hayase (Assistant Professor)
  - ■ Chiemi Watanabe (Assistant Professor)

## ■ Students

- ● Doctoral Program: 8
- ● Master Program: 24
- ● Undergraduate: 7
- ● Research Student: 3

## ■ Adjunct Researchers

- ● Prof. Ishikawa (Nagoya Univ.)
- ● Prof. Ebisawa (JAXA)

# Overview

- R&D in Data Engineering and Databases
- Main Research Areas
  - Information Integration Framework
  - Data Mining and Knowledge Discovery
  - XML and Web Programming
  - Database Applications in Science Domains

# Overview

- **Main Research Areas**
  - **Information Integration Framework**
    - Integration of Heterogeneous Data Sources: DB, Web, File, XML, Sensors, …
    - Stream Processing
      - High-Availability Schemes for Distributed Stream Processing
      - Secure Stream Data Processing
      - Efficient Archiving of Stream Data
      - Outlier Detection over Packet Streams
      - Transactional Stream Processing
    - Indexing for Update-intensive Applications
  - **Data Mining and Knowledge Discovery**
    - Outlier Detection
    - Social Bookmark Analysis
    - Microblog Analysis
    - GPU-based Acceleration of Data Mining

# Overview

- **Main Research Areas (Cont.)**
  - **XML and Web Programming**
    - ・ Online Analytical Processing of XML Data
    - ・ Parallel XML Query Processing using PC-Clusters/Multi-core Processors
    - ・ Faceted-navigation of XML Data
    - ・ Energy-efficient XML Stream Processing
    - ・ RDF/LOD Data Processing
    - ・ Privacy-preserving Database Querying
  - **Database Applications in Science Domains**
    - ・ Development and Maintenance of GPV/JMA Archive
    - ・ Automatic Classification of Pressure Patterns
    - ・ Faceted-Navigation System for QCDml Ensemble XML Data
    - ・ Event Detection from Large Scale Satellite Sensor Data
    - ・ Outburst Detection from X-ray Astronomy Data

# Collaboration

**Japan Aerospace Exploration Agency (JAXA)**

- Outburst Detection from X-ray Astronomy Data

**National Institute of Advanced Industrial Science and Technology (AIST)**

- Event Detection from Large Scale Satellite Sensor Data (GeoGrid)

**Division of Computational Informatics Database Group**

**Division of Global Environmental Science**

- GPV/JMA Data Archive
- Classification of Pressure Patterns

**Division of Particle Physics**

- ILDG/JLDG
- Faceted-Navigation System for QCDml Ensemble XML Data

**Division of Computational Informatics Computational Media Group**

- Regular Division Meeting
- Real-world Sensing Data Management
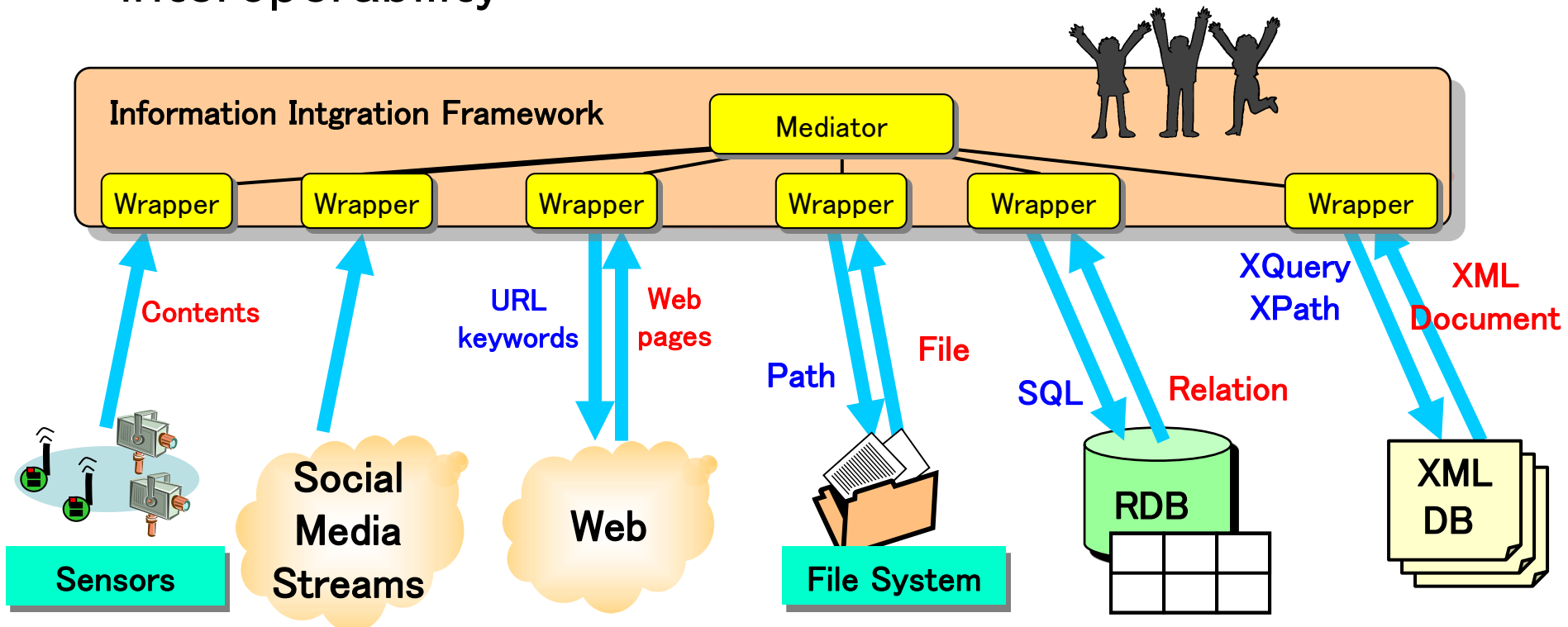
6

# Selected Research Topics

- **Information Integration Framework**
  - Integration of Heterogeneous Data Sources: DB, Web, File, XML, Sensors, …
  - Stream Processing
    - High-Availability Schemes for Distributed Stream Processing
    - ✓ Transactional Stream Processing
    - Secure Stream Data Processing
    - Efficient Archiving of Stream Data
    - Outlier Detection over Packet Streams
  - Indexing for Update-intensive Applications

# Information Integration Framework

- A variety of online data sources
  - Different data formats, access methods, query languages, ⋯
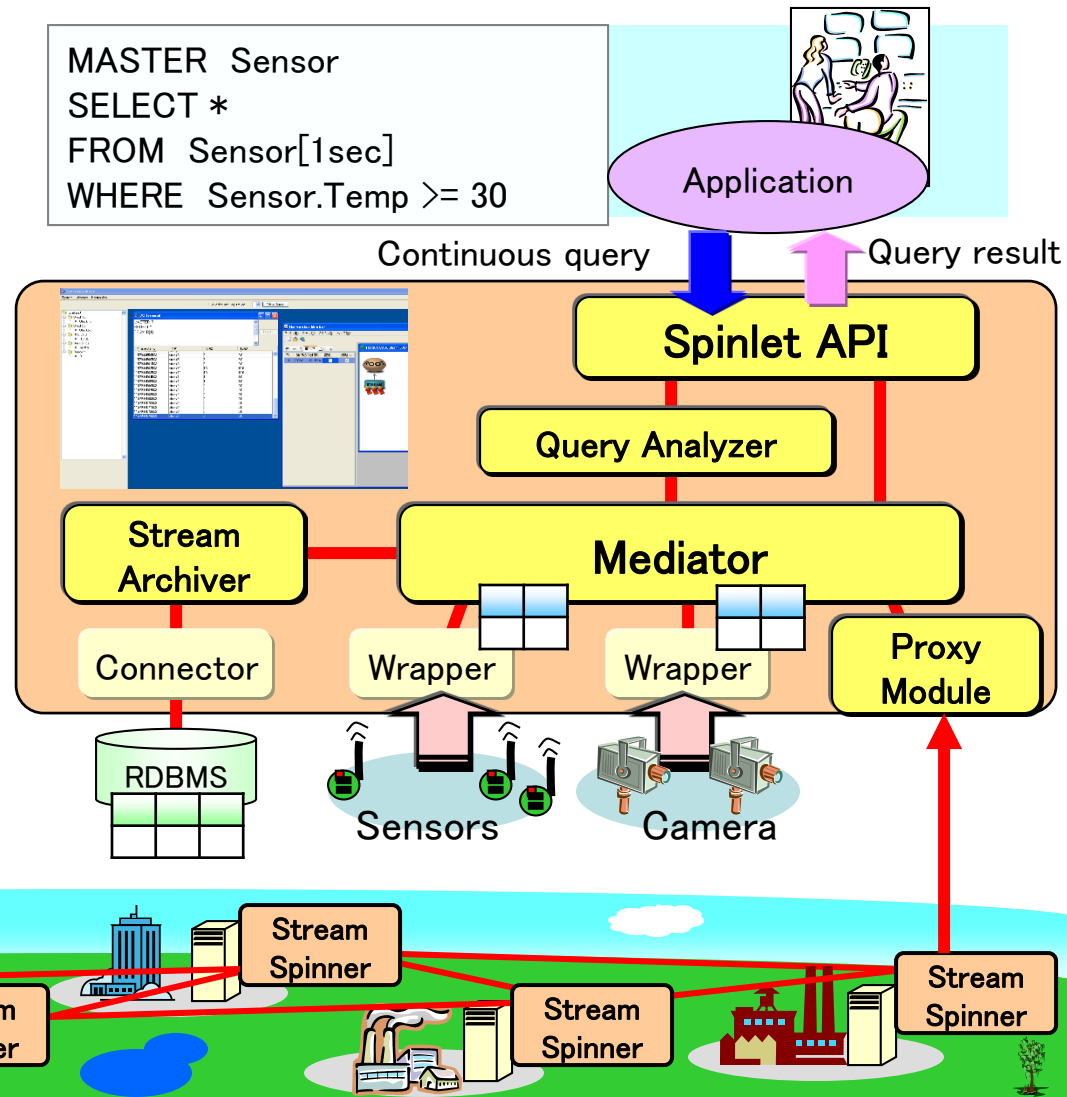- Information integration framework for data interoperability



**Information Intgration Framework**

Mediator

Wrapper | Wrapper | Wrapper | Wrapper | Wrapper | Wrapper

Contents

URL keywords — Web pages

Path — File

SQL — Relation

XQuery XPath — XML Document

Sensors — Social Media Streams — Web — File System — RDB — XML DB

# Data Integration Including Streams

- **StreamSpinner, SS*, JsSpinner**
  - Help integration of heterogeneous data sources
  - Streaming data sources such as sensors, location data, social media streams, etc.
  - Even-driven execution of continuous queries
  - Distributed stream processing



```
MASTER  Sensor
SELECT *
FROM  Sensor[1sec]
WHERE  Sensor.Temp >= 30
```

Application

Continuous query          Query result

Spinlet API

Query Analyzer

Stream Archiver          Mediator

Connector     Wrapper     Wrapper     Proxy Module

RDBMS          Sensors          Camera

Stream Spinner

# Transactional Stream Processing

In data integration context, SPEs do not process only **data streams** but integrates **non-streaming external resources**



Data Stream
Stock Data
Packet

External Resource
Model Data
Databases

SPE

Result Stream

# Problem: Resource reference inconsistency

| External resources may be updated or modified autonomously.

| Continuous query

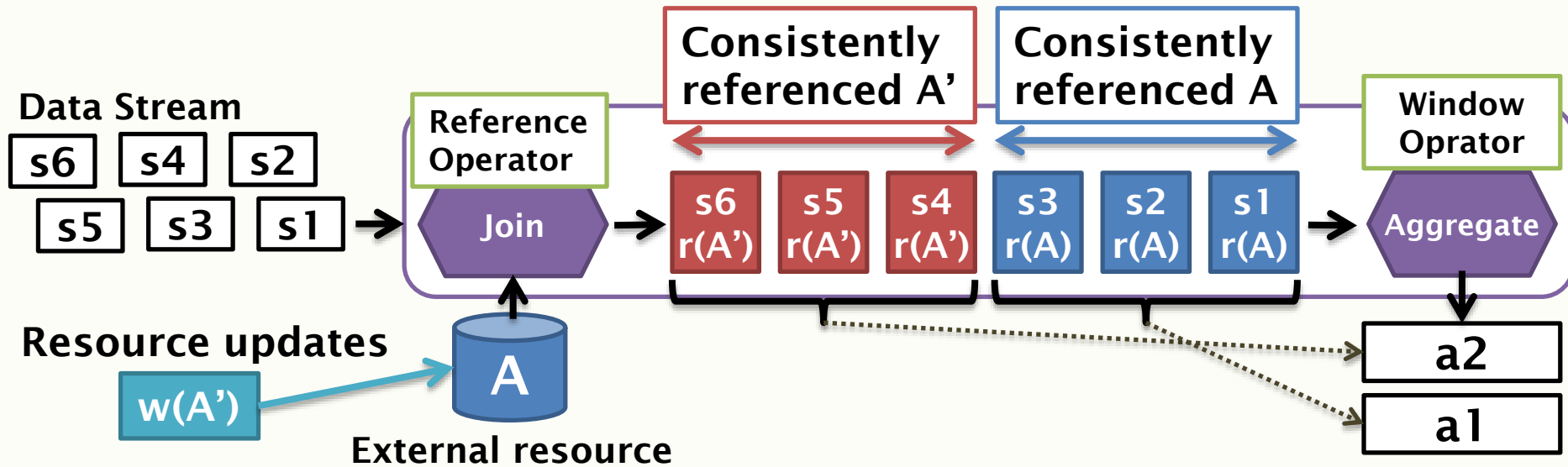  □ **Integrate a data stream & DB** → **aggregate results**
    **(window-based aggregation)**



**Problem**

External resources are not consistently referenced in a single CQ execution instance.

# Our Goal: Transactional Stream Processing

**Goal**

Even if **external resources are updated independently**, we guarantee that **external resources are consistently referenced** in each CQ execution instance.

# How to make stream processing transactional?

## | CQ-processing ensuring serializability
- Serializability of **all CQ-Txns** and DB **Update-Txns**

## | Approaches

(a) Combining exiting concurrency control mechanisms with stream processing

- Two-Phase Locking Strategy (2PL)

- Snapshot Strategy (C2PL)

- Optimistic Strategy

(b) Stream processing combining a redo mechanism and external resource state monitoring

# Selected Research Topics

- **Data Mining and Knowledge Discovery**
  - Outlier Detection
  - Social Bookmark Analysis
  - ✓ Microblog Analysis
  - ✓ GPU-based Acceleration of Data Mining

# Microblog Analysis

■ **Real-world sensing**

  ✓ **Event Detection** [Sakaki+, 10] [Walther+, 13] …
  ✓ **Epidemics Analysis** [Paul+, 11] [Aramaki+, 11] …
  ✓ **Disaster Analysis** [Vieweq+, 10] [Mandel+, 12] …

■ **Location inference**

  - Most users hesitate to disclose their home locations in their profiles.
  - Only few tweets have GEO-tags.

Earthquake!

Shaked!

Shaked!

Um…
Earthquake happened.
But WHERE?

# Graph Based Approach

- Utilize social graphs based on friendships
- Closeness vs. Concentration

**Closeness assumption**
**(Traditional approaches)**

**IF: FRIENDS**
**THEN:  CLOSE**

*Labeled User*

*Unlabeled User*

Tsukuba

*follow*

Tsukuba ?

**Concentration assumption (proposed)**

**IF: FOLLOW A GRAPH LANDMARK**
**THEN:  CLOSE**  A user whose followers are close to each other

*Graph Landmark*

*follow*

Tsukuba

Tsukuba

*Labeled Users*

Tokyo

Tsukuba

*Unlabeled User*

Tsukuba ?

# Graph Based Approach

## Graph Landmark Example



**Boston Fire Dept.**
@BostonFire
Official Twitter Boston Fire. Propane grills on ground or first floor porches only. Charcoal grills on ground only. Never unattended.



red: regular users
blue: graph landmarks

**27%+ IMPROVED**

## Accuracy Comparison

- ✓ LMM:   Proposed  }  concentration assumption
- ✓ UDI:   [Li+, 12]
- ✓ Backstrom: [Backstrom+, 10]
- ✓ Jurgens:  [Jurgens, 13]  }  closeness assumption
- ✓ Naïve:  Medoid

# Content Based Approach
## Static vs. Temporal Local Words

**Static** (traditional)

Local words tightly associated with local regions (city names, home team names, …)

Eagles !

Miyagi

*Labeled Users*

Tigers !

Osaka

*Labeled Users*

Tigers !

*Unlabeled User*

Osaka

**Temporal** (proposed)

Local events (earthquakes, tornados, crimes, …)

Shaked !

*Labeled Users*

Earthquake at Tokyo

Shaked !

*Unlabeled User*

Tokyo if he/she tweets almost at the same time

# Content Based Approach

## Identified Local Events

Distributions of tweets
( a ) at ordinary times
( b ) after an earthquake
    at Hiroshima



( a )        ( b )

**33%+ IMPROVED**

## Accuracy Comparison

✓ Proposed: Proposed        }  temporal
✓ UDI: [Li+, KDD'12]
✓ Cheng: [Cheng+, CIKM'10]   }  static
✓ Random: Random Assignment



Proposed  +
UDI  ×
Cheng  *
Random  □

CDF

Error distance (m)

19

# GPU-based Acceleration of Data Mining

■ Gaining growing attention due to its cost and performance.


Fluid Simulation


Signal Processing


Biology


Database & Data Mining

■ The DB group has been trying to apply GPU for accelerating various data mining processing.

- Probabilistic Latent Semantic Indexing [ICCS'11]
- Frequent itemset mining over uncertain databases [CIKM'12, DEXA'13, IEICE Trans.'14]
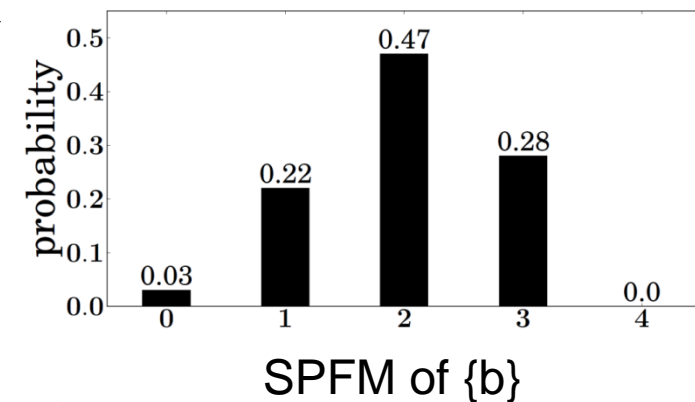- Currently working on sorting huge arrays, time series matching, clustering, and similarity join.

※ Images from NVidia CUDA case home page

# Frequent Itemset Mining over Uncertain Databases

- **A set of transactions**
  - Transaction: ID, itemset, and probability
- **Possible worlds**
  - {T1, T2}: prob. is 0.028
  - {T1, T2, T3}: prob. is 0.252
  - …
- **Support $\mathrm{sup}(X)$**
  - Conventional: # of transactions containing $X$
  - Uncertain: random variable
  - ➜ Support Probability Mass Function (SPMF)
- **Probabilistic Frequent Itemset (PFI)**
  - $P(\mathrm{sup}(X) \geq \mathrm{minsup}) \geq \mathrm{minprob}$
  - minsup and minprob are user-specified values
- **Accelerate frequent itemset mining (PFIM) using GPU**

Uncertain transaction DB

| ID | Itemset | Prob. |
|----|---------|-------|
| T1 | {a, b} | 0.8 |
| T2 | {b, c} | 0.7 |
| T3 | {a} | 0.9 |
| T4 | {a, b, c} | 0.5 |



SPFM of {b}

# SPMF Computation on GPU

$f^1_{\{music\}}$    $f^2_{\{music\}}$    $f^3_{\{music\}}$    $f^4_{\{music\}}$

| 0.2 | 0.8 | 0 | 0 | 0.3 | 0.7 | 0 | 0 | 0.5 | 0.5 | 0 | 0 | 1 | 0 | 0 | 0 |

Convolution

$f^1_{\{music\}} * f^2_{\{music\}}$    $f^3_{\{music\}} * f^4_{\{music\}}$

| 0.06 | 0.38 | 0.56 | 0 | 0.3 | 0.7 | 0 | 0 | 0.5 | 0.5 | 0 | 0 | 1 | 0 | 0 | 0 |

Convolution

$\left(f^1_{\{music\}} * f^2_{\{music\}}\right) * \left(f^3_{\{music\}} * f^4_{\{music\}}\right)$

| 0.03 | 0.22 | 0.47 | 0.28 | 0 | 0 | 0 | 0 | 0.5 | 0.5 | 0 | 0 | 0 | 0 | 0 | 0 |

# Experiments

CPU: Inten Xeon CPU (2.40 GHz) with 4GB memory
GPU: Tesla C2050 (1.15GHz, 3.0GB memory)

■ Accidents
- 106–112x

■ T25I10D500K
- 3.9-22x

# Selected Research Topics

- **XML and Web Programming**
  - Online Analytical Processing of XML Data
  - ✓ Parallel XML Query Processing using PC-Clusters/Multi-core Processors
  - Faceted-navigation of XML Data
  - Energy-efficient XML Stream Processing
  - ✓ RDF/LOD Data Processing
  - Privacy-preserving Database Querying

# Parallel XML Query Processing on a Multi-core System

- Pattern matching queries are important in querying XML.
- Holistic twig joins (HTJ)
  - A family of XML query processing algorithms
  - Find matches for a given query tree (twig)
- Propose a parallel version of TwigStack algorithm for multi-core processors

# TwigStack Algorithm



**XML Database**

**XML Query**

```
        name
       /    \\
   lname    fname
     |        |
   kita     nishi
```

**XML Node Streams**

| n1 | n2 | n3 | n4 | n5 |
|----|----|----|----|----|

| ln1 | ln2 | ln3 | ln5 |
|-----|-----|-----|-----|

| kt1 | kt2 | kt4 |
|-----|-----|-----|

| fn1 | fn3 | fn4 | fn5 |
|-----|-----|-----|-----|

| ns1 | ns3 |
|-----|-----|

**1st Phase (Task 1):**
1. Perform root-to-leaf pattern matching.
2. Generate root-to-leaf path solutions.

√ n1-ln1-kt1
√ n1-fn1-ns1
  n2-ln2-kt2
  n3-fn3-ns3

**2nd Phase (Task 2):**
1. Compile root-to-leaf path solutions for final solutions.

n1-ln1-kt1-fn1-ns1

Tree structure:

club (1, 1:55, 1)
- clubname (1, 2:4, 2)
  - soccer (1, 3, 3)
- member (1, 5:14, 2)
  - name n1
    - lname ln1
      - kita kt1
    - fname fn1
      - nishi ns1
- member (1, 15:24, 2)
  - name n2
    - lname ln2
      - kita kt2
- member (1, 25:34, 2)
  - name n3
    - lname ln3
      - minami (1, 28, 5)
    - fname fn3
      - nishi ns3
- member (1, 35:44, 2)
  - name n4
    - fname fn4
      - kita kt4
- member (1, 45:54, 2)
  - name n5
    - lname ln5
      - higashi (1, 48, 5)
    - fname fn5

# Basic Idea

- Partition the XML tree and process in parallel.



**How to partition XML node streams so that each core can derive path solutions independently.**

# Experiments

- **Sequential execution time**
  - Q1: 18.48 s, Q2: 44.62 s, Q3: 17.16 s, Q4: 24.65 s, Q5: 15.44 s

## 4 GB of XML Data



Simplest Structure — Lowest Selectivity

□ Data Parallelism  ■ Task Parallelism          **#CPU-Cores**

# Linked Open Data, RDF

- **Linked Open Data (LOD) is increasing rapidly**
  - A method to publish and share structured data on the Web
  - "Web of Data": Data linked with each other
- **Resource Description Framework (RDF)**
  - A framework for describing resources on the Web
  - Triple: Subject, Predicate, and Object
- **Numerical data also published as Linked Open Data**
  - Statistics from governments, sensor data, etc.
  - Growing demands for analytical processing over LOD data..
- **We propose an ETL framework for the OLAP analysis of LOD datasets.**
  - Derivation of a star schema from a large RDF graph.

# Framework Overview

- Generated schema from LOD dataset.
  - fact-table) observation_instance
  - dim.-table) Time, Location

Environmental radio activity level monitoring data published by Japan Nuclear Regulation Authority

**Obs. Instance**

ra:20110315/p02/t20 — rdf:value → "0.040" — **Value**

ev:time

**time**

time:20110315T22PT1H — tl:at → "2011-04-14T00:00:00" ^^xsd:dateTime

ev:place

**location**

gn:2111833

**Dim.**

| place | |
| --- | --- |
| subject | |
| layer_1 （district） | |
| layer_2 （prefecture） | |
| layer_3 （country） | |
| layer_4 （continent） | |

**Fact**

| observation_instance |
| --- |
| subject |
| place_subject |
| time_subject |
| value |

← **Measure**

**Dim.**

| time |
| --- |
| subject |
| sec |
| min |
| hour |
| day |
| month |
| year |

# Selected Research Topics

- **Database Applications in Science Domains**
  - Development and Maintenance of GPV/JMA Archive
  - Automatic Classification of Pressure Patterns
  - ✓ Faceted-Navigation System for QCDml Ensemble XML Data
  - Event Detection from Large Scale Satellite Sensor Data
  - ✓ X-ray Outburst Detection from X-ray Astronomy Data

# Int'l Lattice Data Grid (ILDG)

■ An international collaboration which provides standards, services, methods and tools that facilitates the sharing and interchange of lattice QCD gauge configurations by integrating their regional data grids.

■ File formats in ILDG
  ● Configuration binary (10+TB in JLDG)
    ・ LIME (Lattice QCD Interchange Message Encapsulation)
  ● Metadata (QCDml)
    ・ Ensemble XML (200+ in ILDG)
    ・ Configuration XML (30,000+ in JLDG)
  ● A number of configuration binaries are associated with an ensemble in terms of markovChainURI and LFN.

# QCDml Ensemble XML

```xml
<markovChain xmlns="…">
  <markovChainURI>mc: //JLDG/CP-PACS/RCNF2/RC12x24-
B1800K014090C1600</markovChainURI>
  <management>
    <revisions>1</revisions>
    <collaboration>CP-PACS</collaboration>
    <projectName>RCNF2 (Nf=2 full QCD with iwasaki RG gauge and
tadpole improved clover quark action)</projectName>
    <ensembleLabel>B1800</ensembleLabel>
    <reference>Phys.Rev. D65 (2002) 054505 (hep-lat/0105015),
Erratum-ibid. D67 (2003) 059901</reference>
    <archiveHistory>
      <elem>
        <revision>1</revision>
        <revisionAction>add</revisionAction>
        <participant>
          <name>T.Yoshie</name>
          <institution>Center fof Computational Sciences, University
of Tsukuba</institution>
```

# QCDml Faceted Navigation Interface System Overview



ILDG

Facet Navigation System
(PHP + SQL + XQuery)

Web Server
(Apache)

JLDG

USQCD

LDG

UKQCD

CSSM

Facet Database

Facet extraction (XQuery)

QCDml
Ensemble (ILDG)
& Configuration (JLDG)

Downloading Ensemble XML

RDBMS (MySQL)

XML DB (eXist)

# Demonstration

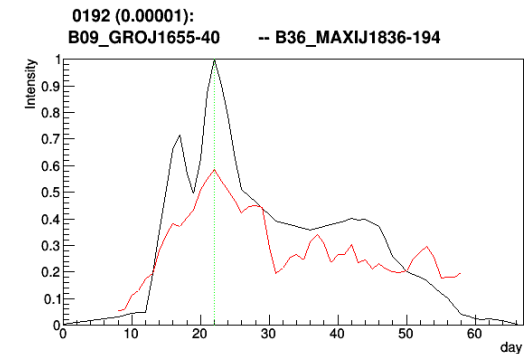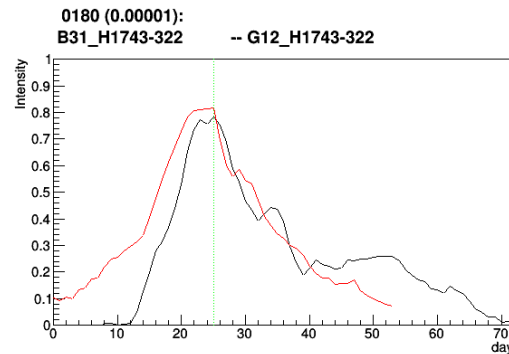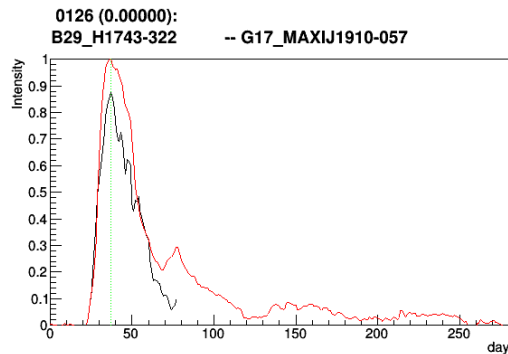# Similarity Search over Light Curves of X-ray Outbursts

- A collaborative work with JAXA.

- X-ray outbursts
  - Phenomena in which X-ray emission from a celestial object grows for a certain time period.

- Researchers in JAXA are interested in finding celestial objects showing similar light curves.

GX339-4, XTEJ1752-223

# Our Work

- We apply time-series analysis techniques, such as DTW and DDTW, to search for similar light curves out of massive observation data.

- Detected light curves:

# Major Funding

- Grant-in-Aid for Scientific Research from Ministry of Education, Culture, Sports, Science and Technology (MEXT) (~$1 Million; past 6 years)
  - Grant-in-Aid for Scientific Research A
  - Grant-in-Aid for Scientific Research on Priority Areas (Infoplosion Project)
  - Grant-in-Aid for Exploratory Research
  - Grant-in-Aid for Young Scientists
- MEXT Big Data Federation Feasibility Study (~$0.3 Miillion; 2013)
- From industry

# Collaboration

- **Industries**
  - Hitachi
  - NEC
  - Fujitsu Lab.
  - Mitsubishi Electric
  - NTT Lab.
  - KDDI Lab.

  etc.
- **International**
  - Carnegie Mellon University
  - Chinese Univ. of Hong Kong
  - Georgia Institute of Technology

  etc.

# Publication and Awards

- **Refereed Papers**
  - 2008: 29 (Journal 8, Conference 21 (3 Demo/Posters))
  - 2009: 28 (Journal 14, Conference 14 (1))
  - 2010: 22 (Journal 7, Conference 15)
  - 2011: 16 (Journal 4, Conference 12 (3))
  - 2012: 21 (Journal 8, Conference 13 (3))
  - 2013: 18 (Journal 4, Conference 14 (4))
- **Awards**
  - 4 Best Paper Awards (IEICE Trans., iiWAS2010, IPSJ SIG, DBSJ Journal)
  - 4 Best Student Paper Awards (WAIM2008, iiWAS2008, KMIS2010)
  - Contribution Award (IEICE)
  - 26 Students' Awards

# Future Plan

- ■ Research and Development for Data Engineering Challenges
  - ● Data integration framework to accommodate Big Data.
  - ● Big Data analysis challenges.
  - ● New issues involving social media and open data: privacy, social readings, LOD (Linked Open Data).
- ■ Database Applications in Science Domains
  - ● Started collaboration with Biological Science Group through the dual degree program on gene databases
- ■ Reinforcement of cooperation with other divisions and organizations

# Thank you.