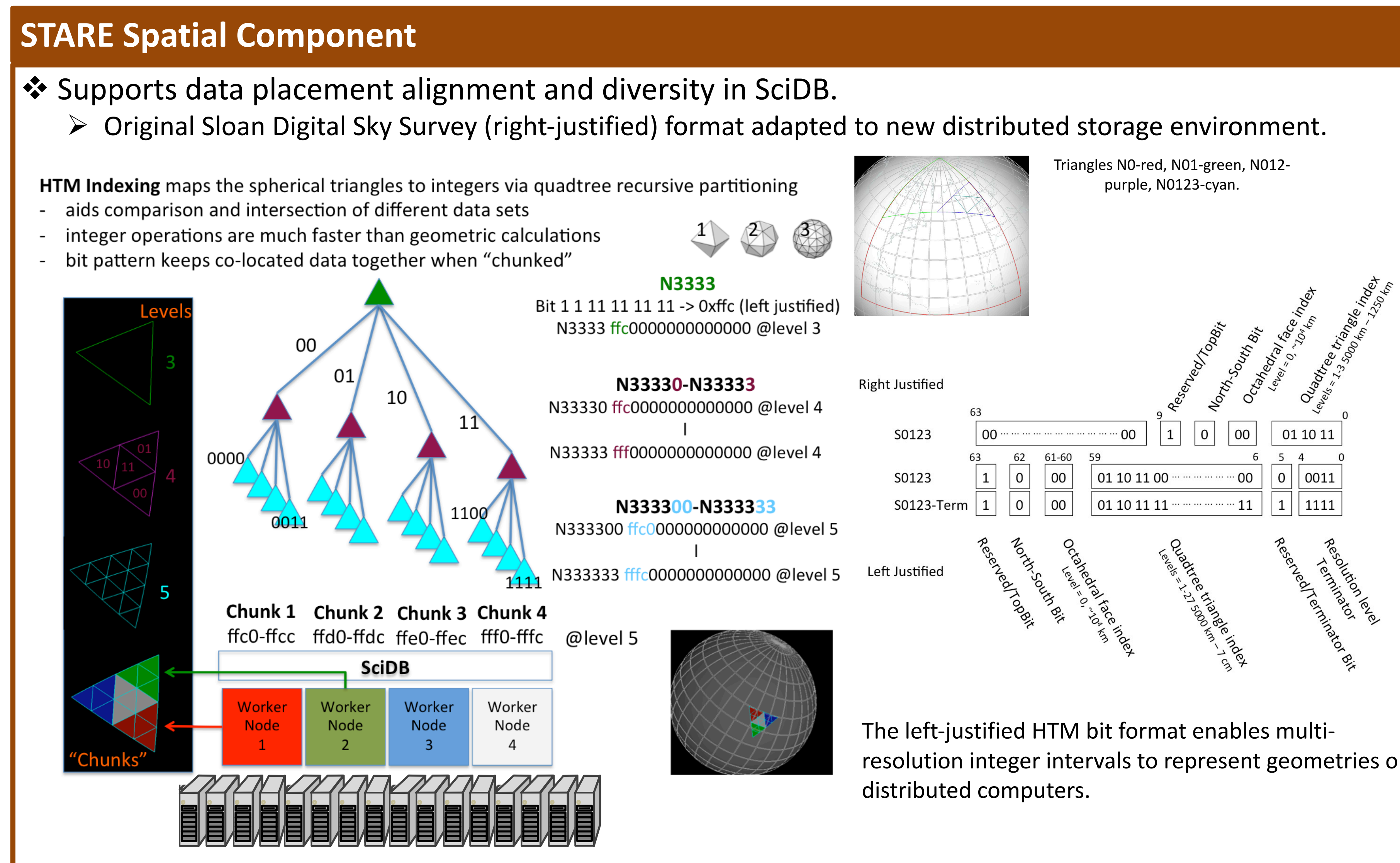# An Innovative Infrastructure with a Universal Geo-spatiotemporal Data Representation Supporting Cost-effective Integration of Diverse Earth Science Data

Michael Lee Rilee[1,3] and Kwo-Sen Kuo[2,3]    [1]Rilee Systems Technologies LLC; [2]Bayesics LLC; [3]NASA Goddard Space Flight Center

## Abstract

❖ The SpatioTemporal Adaptive Resolution Encoding (STARE) is **a unifying scheme encoding geospatial and temporal information** for organizing data on scalable computing/storage resources, minimizing expensive data transfers.

❖ STARE provides a compact representation that **turns set-logic functions into integer operations**, e.g. conditional subsetting, taking into account representative spatiotemporal resolutions of the data in the datasets.

❖ STARE **geo-spatiotemporally aligns data placements of diverse data** on massive parallel resources to maximize performance.

❖ **Automating important scientific functions (e.g. regridding) and computational functions (e.g. data placement)** allows scientists to focus on domain specific questions instead of expending their efforts and expertise on data processing.

❖ With STARE-enabled automation, SciDB+STARE provides a database interface, reducing costly data preparation, increasing the volume and variety of interoperable data, and easing result sharing.

❖ Using SciDB+STARE as part of an integrated analysis infrastructure dramatically eases combining diametrically different datasets.

## STARE Spatial Component

➤ Supports data placement alignment and diversity in SciDB.
  ➤ Original Sloan Digital Sky Survey (right-justified) format adapted to new distributed storage environment.

**HTM Indexing** maps the spherical triangles to integers via quadtree recursive partitioning
- aids comparison and intersection of different data sets
- integer operations are much faster than geometric calculations
- bit pattern keeps co-located data together when "chunked"

Triangles N0-red, N01-green, N012-purple, N0123-cyan.

**N3333**
Bit 1 1 11 11 11 11 -> 0xffc (left justified)
N3333 ffc0000000000000 @level 3

**N33330-N33333**
N33330 ffc0000000000000 @level 4
|
N33333 fff0000000000000 @level 4

**N333300-N333333**
N333300 ffc0000000000000 @level 5
|
N333333 fffc000000000000 @level 5

Chunk 1    Chunk 2    Chunk 3    Chunk 4
ffc0-ffcc  ffd0-ffdc  ffe0-ffec  fff0-fffc   @level 5

SciDB

| Worker Node 1 | Worker Node 2 | Worker Node 3 | Worker Node 4 |

"Chunks"

**Right Justified**
S0123 | 00 ... 00 | 1 | 0 | 01 10 11 |
S0123 | 1 | 0 | 00 | 01 10 11 00 | 00 | 0 | 0011 |
S0123-Term | 1 | 0 | 01 10 11 11 | 11 | 1 | 1111 |

Reserved/right — North-South bit — Octahedral face index (Levels 0, 1.25e4 km) — Quadtree triangle index (Levels 1 - 27 2.5e4 km - 3 cm) — Reserved/terminator bit — Resolution level

**Left Justified**

The left-justified HTM bit format enables multi-resolution integer intervals to represent geometries on distributed computers.

## STARE Temporal Component

❖ Hierarchical like HTM
❖ Bit ranges capture important time resolutions
❖ Implementation flexible, supporting multiple schema for bit ranges
❖ Not (necessarily) a tree, but
❖ Supports intervals like HTM

Meanings of bit ranges
(from least to most significant)
*prototype*

| Range | Starting Bit | Ending Bit | No. Bits | Denoting |
|---|---|---|---|---|
| 0 | 0 | 2 | 3 | *Resolution* |
| 1 | 3 | 12 | 10 | millisecond |
| 2 | 13 | 24 | 12 | Second |
| 3 | 25 | 29 | 5 | Hour |
| 4 | 30 | 32 | 3 | Day of week |
| 5 | 33 | 34 | 2 | Week |
| 6 | 35 | 38 | 4 | Month |
| 7 | 39 | 48 | 10 | Year |
| 8 | 49 | 58 | 10 | Kilo-annum |
| 9 | 59 | 62 | 4 | Mega-annum |
| 10 | 63 | 63 | 1 | Before/After |

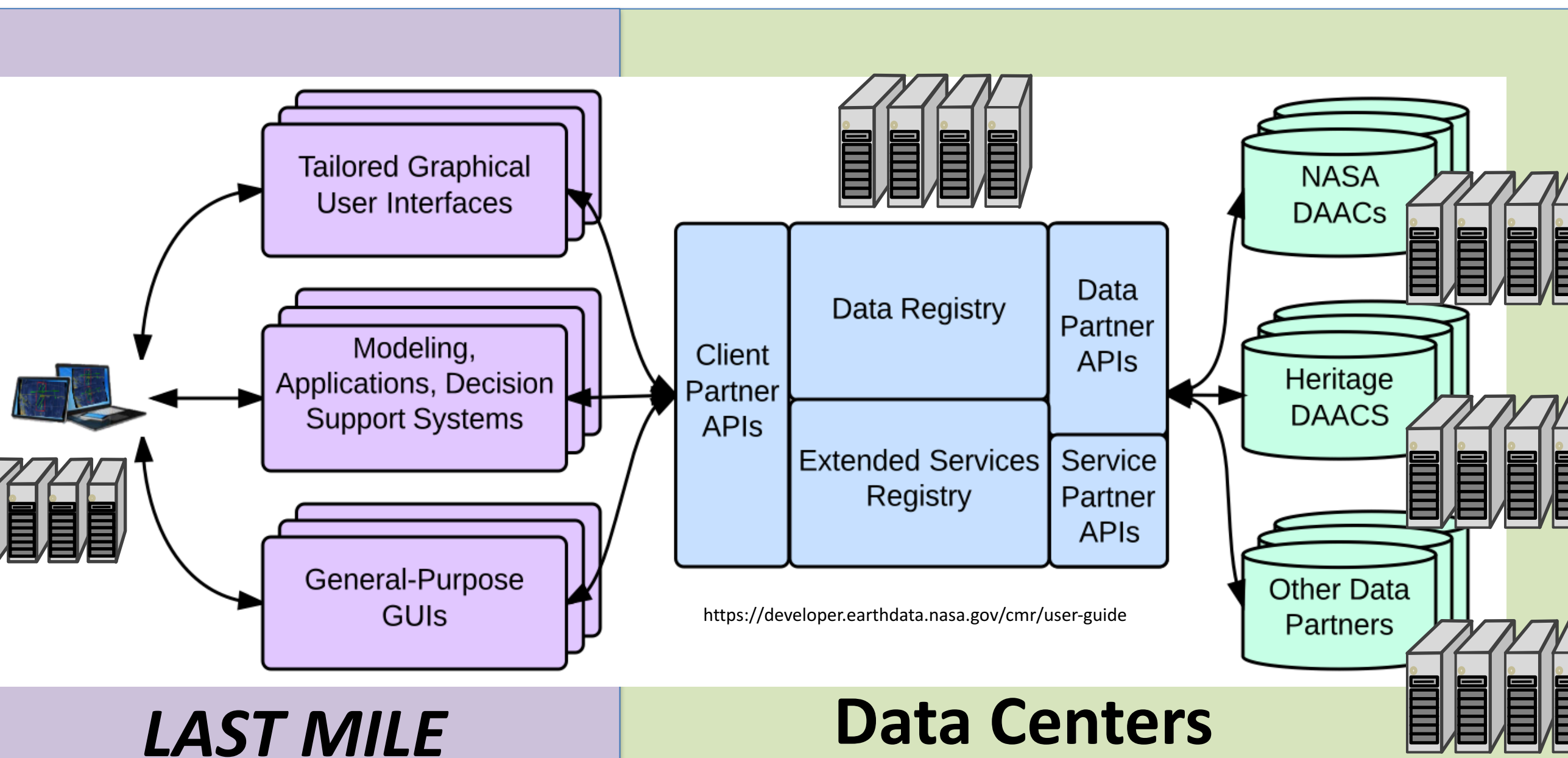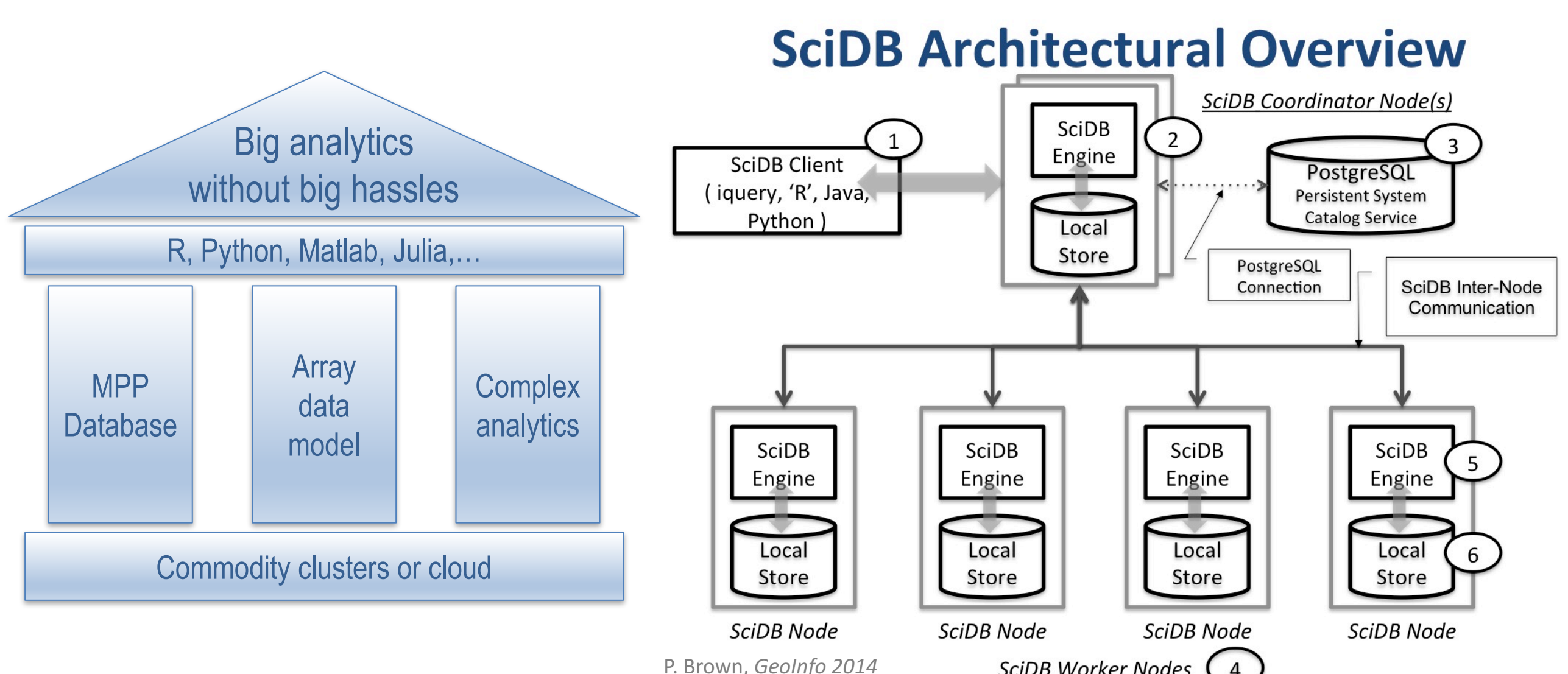## STARE example metadata for a MODIS granule

```
<notional-modis-stare-metadata>
<TemporalHex> x404e0052c0003
</TemporalHex>
<SpatialHex>
x4a00000000000004,
x4a20000000000004 x4a3fffffffffff,
x4a50000000000004,
x4a80000000000004 x4adfffffffffff,
x4af0000000000004,
x4b10000000000004,
x4b60000000000004,
x4b0e0000000000004,
x4c00000000000004 x4c1fffffffffff,
x4c30000000000004 x4c7fffffffffff,
x4ca0000000000004,
x4cc0000000000004,
x4ce0000000000004 x4cffffffffffff,
x4f80000000000004,
x4fa0000000000004 x4fbfffffffffff
</SpatialHex>
</notional-modis-stare-metadata>
```

*Resolution level reduced for clarity*

## Why SciDB?

### Resource Consumption Advantages
❖ Minimize download and local data management
❖ Free end-user resources for research and science

### Performance Advantages
❖ Array data model is **better suited for scientific data than relational databases**.
❖ Tightly coupled analysis and storage layers allows **better optimization than Spark**.

### SciDB Architectural Overview

Big analytics without big hassles
R, Python, Matlab, Julia,…

| MPP Database | Array data model | Complex analytics |

Commodity clusters or cloud

SciDB Client (iquery, 'R', Java, Python) ① → SciDB Engine / Local Store — *SciDB Coordinator Node(s)* ③ PostgreSQL Persistent System Catalog Service
PostgreSQL Connection — SciDB Inter-Node Communication

SciDB Engine / Local Store (×4) — *SciDB Worker Nodes* ④ ⑤ ⑥

SciDB Node    SciDB Node    SciDB Node    SciDB Node

P. Brown, GeoInfo 2014

**LAST MILE**

Tailored Graphical User Interfaces
Modeling, Applications, Decision Support Systems
General-Purpose GUIs

Client Partner APIs
Data Registry
Data Partner APIs
Extended Services Registry
Service Partner APIs

**Data Centers**
NASA DAACs
Heritage DAACS
Other Data Partners

https://developer.earthdata.nasa.gov/cmr/user-guide

## SciDB Query

**To spatiotemporally "join" the two datasets with STARE indexing (5-min at level 7, i.e. ~78-km resolution):**
    join( nmq_precip, trmm1_2B31 ); - **Magic!**

**To subset temporally for visualization:**
    select *
    into nmq_trmm_09120303
    from nmq_trmm1_result - result from previous join
    where
      tIndex= temporalIndexFromString("2009-11-03 03:00:00.000 (00)") or
      tIndex= temporalIndexFromString("2009-11-03 03:04:16.000 (00)");

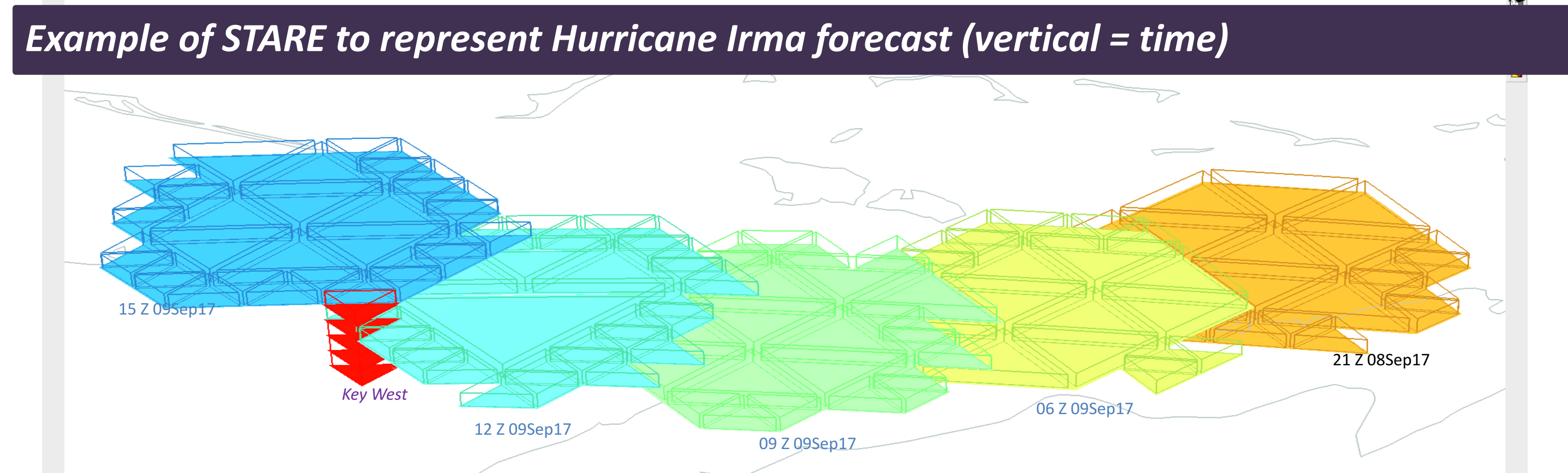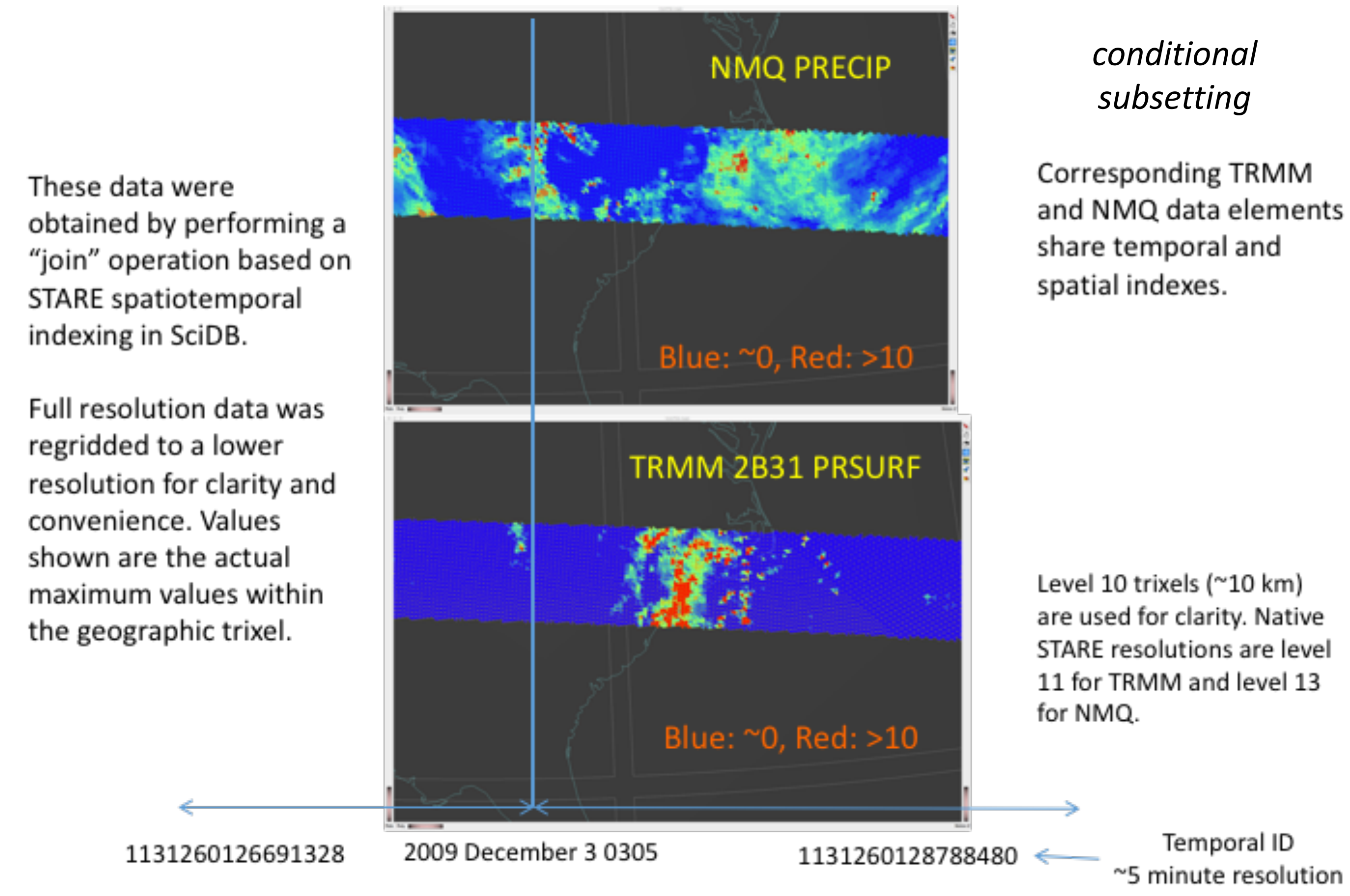**Joining a week of NMQ and TRMM 2B31 takes less than 1 minute on MAS cluster.**
    NMQ: 2016(=12*24*7) 5-min time slices of 7000x3500 2D floating-point array.
    TRMM 2B31: 110 orbit granules of ~9000x49 2D floating-point array (multiple attributes)
No more painstakingly reading and filtering numerous files from various datasets!
    Operations can be straightforwardly interfaced with a GUI for *Visual Analytics*!
**One remaining performance hurdle is *visualization* – data movement! (see summary for related presentations)**

These data were obtained by performing a "join" operation based on STARE spatiotemporal indexing in SciDB.

Full resolution data was regridded to a lower resolution for clarity and convenience. Values shown are the actual maximum values within the geographic trixel.

**NMQ PRECIP**
Blue: ~0, Red: >10

**TRMM 2B31 PRSURF**
Blue: ~0, Red: >10

*conditional subsetting*

Corresponding TRMM and NMQ data elements share temporal and spatial indexes.

Level 10 trixels (~10 km) are used for clarity. Native STARE resolutions are level 11 for TRMM and level 13 for NMQ.

1131260126691328    2009 December 3 0305    1131260128788480
Temporal ID ~5 minute resolution

## Example of STARE to represent Hurricane Irma forecast (vertical = time)

15 Z 09Sep17    12 Z 09Sep17    09 Z 09Sep17    06 Z 09Sep17    21 Z 08Sep17
Key West

❖ **STARE for spatiotemporal regions for data search and combination**
  ➤ Multi-resolution shows flexibility of the scheme
  ➤ Geodesic edges speed geometric calculations
  ➤ Very general regions and temporal structure can be supported
  ➤ Useful for metadata and for general spatiotemporal specification
  ➤ Memory and compute efficient
  ➤ Naturally supports efficient data placement and parallel computing
  ➤ Supports processing closer to where data are stored

## Summary

### SciDB on STARE
❖ Provides a unifying scheme for comparing and combining diverse data sets
❖ Naturally supports data placement alignment for efficient use of SciDB
❖ Set and logic operations are efficient, straightforward to code
❖ Transparent use of high-end parallel/distributed compute & storage
❖ Scientists can work with data via high-level queries
❖ Growing set of functions for representation, regridding in SciDB enabled by STARE

### Related Presentations:
STARE in Visualization: IN23F-07, IN33C-0141; in the SciDB array database: IN41B-0035 (this work), IN33E-04, and the path to enabling machine learning: IN11E-07.