



<u>http://l-sis.org</u> → publications

http://en.wikipedia.org/wiki/Array_DBMS



Who Needs Arrays?

- Sensor, image, simulation, statistics data
 - Earth: Geodesy, geology, hydrology, oceanography, climate, earth system, ...
 - Space: optical / radio astronomy, cosmological simulation, planetary science, ...
 - Life: Pharma/chem, healthcare / bio research, bio statistics, genetics, ...
 - Engineering & research: Simulation & experimental data in automotive/shipbuilding/ aerospace industry, turbines, process industry, ...
 - **Management/Controlling:** Decision Support, OLAP, Data Warehousing, census, statistics in industry and public administration, ...
 - Multimedia: distance learning, prepress, ...
- "80% of all data have some spatial connotation" [C&P Hane, 1992]





CONCEPTUAL MODELLING



Array Analytics

Array Analytics :=

Efficient analysis on multi-dimensional arrays of a size several orders of magnitude above evaluation engine's main memory

- Essential data property: n-dimensional Cartesian neighborhood
 - Secondary: #dimensions, density, ...
- Operations: signal/image processing, Linear Algebra [M. Stonebraker], iterations





The Array Data Model





SYSTEMS

Early History of Array Databases



Array DBMSs Landscape Today

rapidly evolving ecosystem → necessarily incomplete

- Array Database Systems
 - query language, multi-user operation, storage management, access control
 - Ex: rasdaman, SciDB, EXTASCID, PostGIS Raster, Oracle GeoRaster
- Array tools: command-line tools & libraries, but no service
 - no query concept, but procedural API
 - Ex: OpenDataCube, OPeNDAP, Wendelin.core, TensorFlow, boost::geometry, xtensor, TileDB, ArrayStore, Ophidia
- Map/Reduce: Hadoop & Spark as cloud parallelization paradigm
 - Array layers on top of Hadoop, Spark
 - Ex: SciHadoop, Spatial Hadoop, GeoTrellis, MrGeo, SciSpark, ClimateSpark



Array DBMSs Landscape Today

- Array Database Systems
- Array tools: command-line tools & libraries, but no service
- Map/Reduce: Hadoop & Spark as cloud parallelization paradigm

- <u>Technology overview</u>
 - 19 technologies compared
 - 4 benchmarked



rasdaman

- "<u>ras</u>ter <u>da</u>ta <u>man</u>ager": SQL + n-D arrays
 - Scalable parallel "tile streaming" architecture
 - [VLDB 1994, VLDB 1997, SIGMOD 1998, VLDB 2003, ..., VLDB 2016]
- Blueprint for stds, in operational use





The rasql Query Language

selection & subsetting

select c.data[*:*, 100:200, *:*, 42]
from ClimateSimulations as c

- result processing
 select img.data * (img.data.green > 130)
 from LandsatArchive as img
- search & aggregation



data format conversion

select encode(c.data[*:*,*:*,100,42], "png")
from ClimateSimulations as c





rasdaman

DB

 \rightarrow





C>ONSTRUCTOR

Linear Algebra Ops



Histogram

select marray bucket in [0:255]
 values count_cells(img = bucket)
from img



[SSDBM 2014]

Arrays in SQL



rasdaman as blueprint

scene: row(band1: integer, ..., band7: integer) mdarray [0:4999,0:4999])



ARCHITECTURE



Storage Management



- Divergent access patterns for ingest and retrieval
- Server must mediate between access patterns

Adaptive Partitioning ("Tiling")

- Any tiling [Furtado 1999]
 - Cast into strategies

 rasdaman storage layout language

Why irregular tiling?



insert into MyCollection
values ...
tiling
area of interest [0:20,0:40], [45:80,80:85]
tile size 1000000
index d_index storage array compression zlib



Query Processing

- Clear separation: set vs array trees
 - Arrays as 2nd order attributes
- Optimization
- Tile-based evaluation



select a.array < sum cells(</pre>

b.array + c.array)

a, b, c

from

C>ONSTRUCTOR

UNIVERSITY



Query Optimization

select max_cells(a + b)
from a, b



[Ritsch 2000]



Parallel / Distributed Query Processing

select

max((A.nir - A.red) / (A.nir + A.red))- max((B.nir - B.red) / (B.nir + B.red)) - max((C.nir - C.red) / (C.nir + C.red)) from A, B, C Dataset C 1 query \rightarrow 1,000+ cloud nodes [ACM SIGMOD DANAC 2014] **Dataset A** Dataset B

C>ONSTRUCTOR

UNIVERSITY



APPLICATIONS

Gene Expression Analysis

http://urchin.spbcas.ru/Mooshka/ [Samsonova et al]

- Gene expression = reading out genes for reproduction
- Research goal: capture spatio-temporal expression patterns in Drosophila



select encode(scale({1c,0c,0c}*e[0,*:*,*:*]
 +{0c,1c,0c}*e[1,*:*,*:*]
 +{0c,0c,1c}*e[2,*:*,*:*], 0.2), "image/jpeg")
from EmbryoImages as e
where oid(e)=193537

Human Brain Imaging

- Research goal: structural-functional relations in human brain
- Experiments \rightarrow activity patterns (PET, fMRI)
 - Temperature, electrical, oxygen consumption, ...
 - → lots of computations → "activation maps"
- Example: "a parasagittal view of all scans containing critical Hippocampus activations, TIFF-coded."

select tiff(ht[\$1, *:*, *:*])
from HeadTomograms as ht,
Hippocampus as mask
where count_cells(ht > \$2 and mask)
 / count_cells(mask)
 > \$3

\$1 = slicing position, \$2 = intensity threshold value, \$3 = confidence



C>ONSTRUCTOR

UNIVERSITY



Cosmological Simulation

- Modelling domain: 4D
- Results: 3D/4D cutouts from universe
- Screenshots: AstroMD [Gheller, Rossi 2001]



Early 3-D Service on rasdaman

[Diedrich et al 2001]







- Agile Analytics on x/y/t + x/y/z/t Earth & Planetary datacubes
 - EU rasdaman
 + US NASA WorldWind
 - Rigorously standards as c/s APIs

earthserver.world

- Multi-Petabyte, worldwide
 - participation free & open





Co-funded by the European Union



DEMO

Edge Integration & Fusion

- Airborne drone
 - Demonstrated at NATO C-UAS Exercise
 - Relatime downlink
- Ship
 - Demonstrated on research vessel, Nuuk/Greenland
 - Realtime send & receive
- Nanosat
 - Demonstrated in orbit
 - Faster: avoid full download
 - QoS: Deliver answers, not pixels



AI + Datacube Integration

- Seamless integration of ML in datacube queries
 - Ex: model-based crop classification, fully integrated

```
for $s2 in (Sentinel_2),
    $m in (CropModel)
return encode( nn.predict( $s2[...], $m ), "tiff" )
```

• Natural Language Processing (RSVQA, Begüm Demir)



C>ONSTRUCTOR

UNIVERSITY

 WCPS Chatbot for geo datacube analytics

Say Hello to Datacubes

Let AI write & explain datacube analytics



WRAP-UP



Summary

- Arrays are core data structure next to sets, graphs, hierarchies
 - sensor, image, simulation, statistics datacubes
- Array DBMS for declarative queries on massive n-D arrays
 - rasdaman
- Issues:
 - enhancing distributed processing
 - iterative methods





Advertisement

- Seeking datacube coders
 - Thesis see my group's <u>current list of thesis topics</u> Sabbatical in Spring 2025
 - Research projects
- Common requirement: strong coding skills
 - JavaScript / TypeScript / frameworks; Java; C++