

# The Era of Big Spatial Data

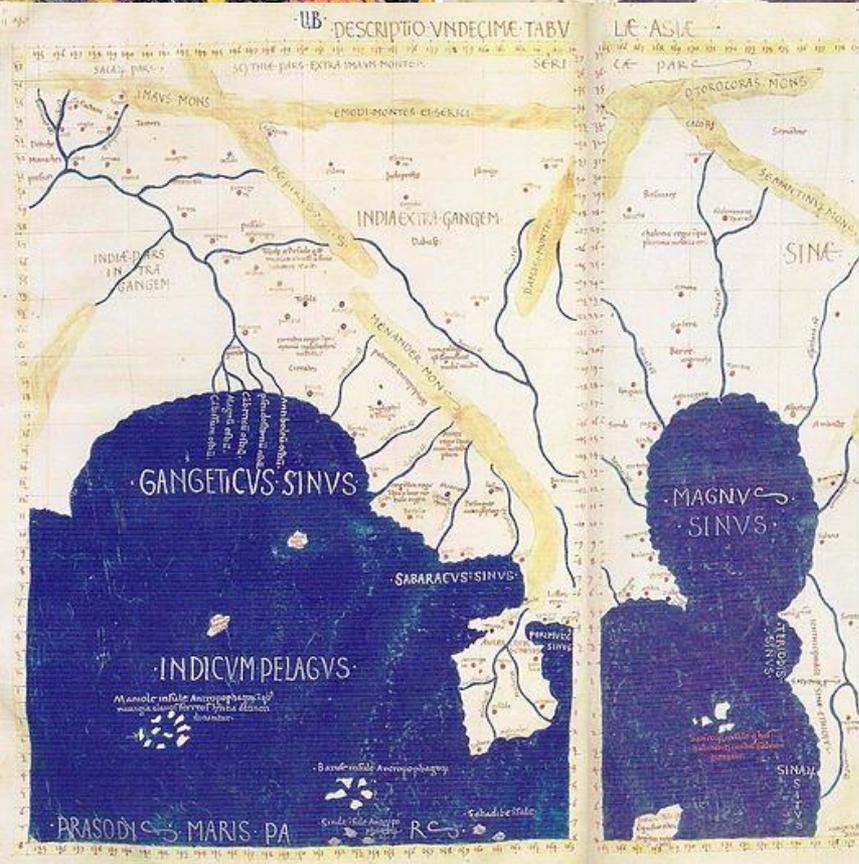
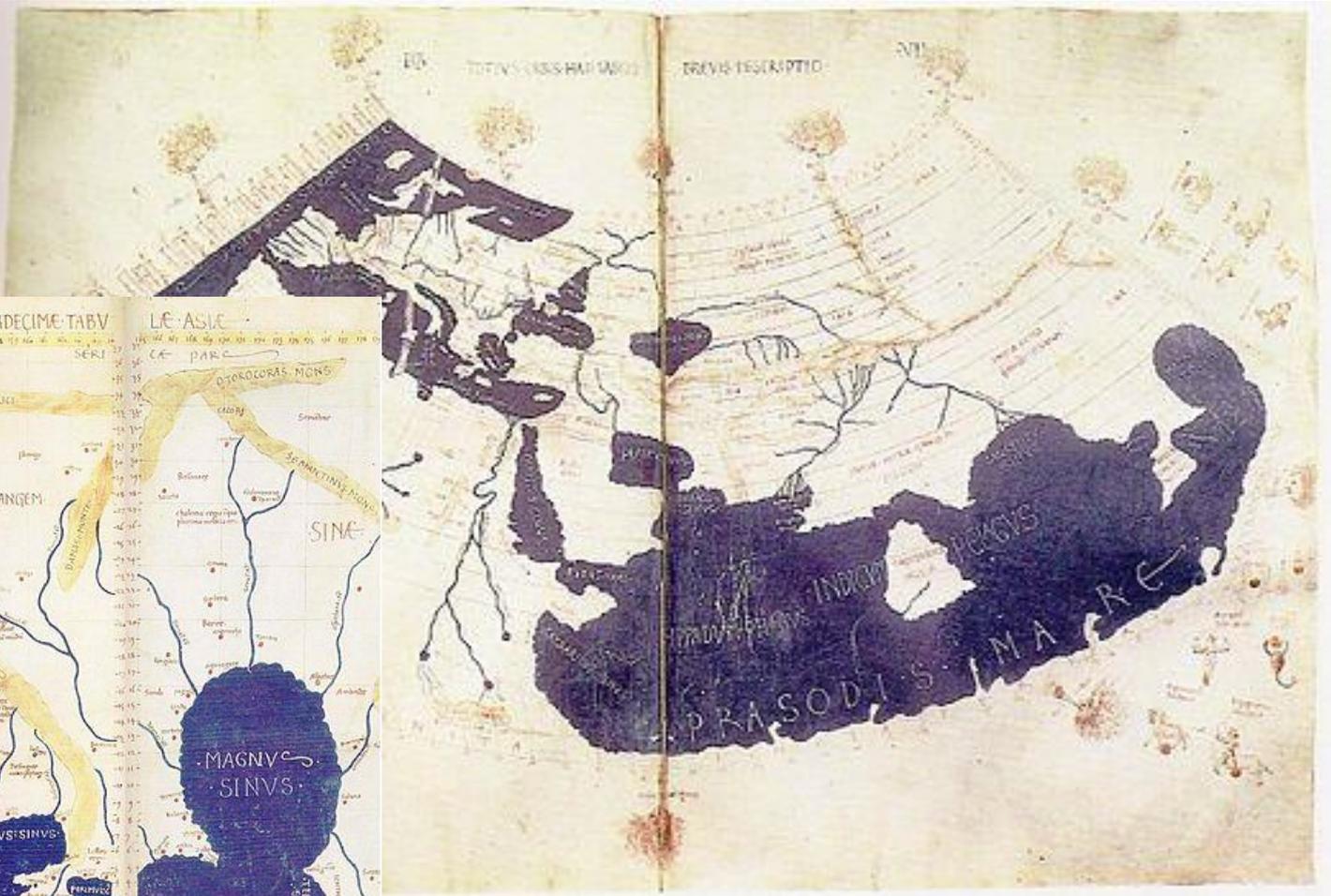
Ahmed Eldawy

Computer Science and Engineering

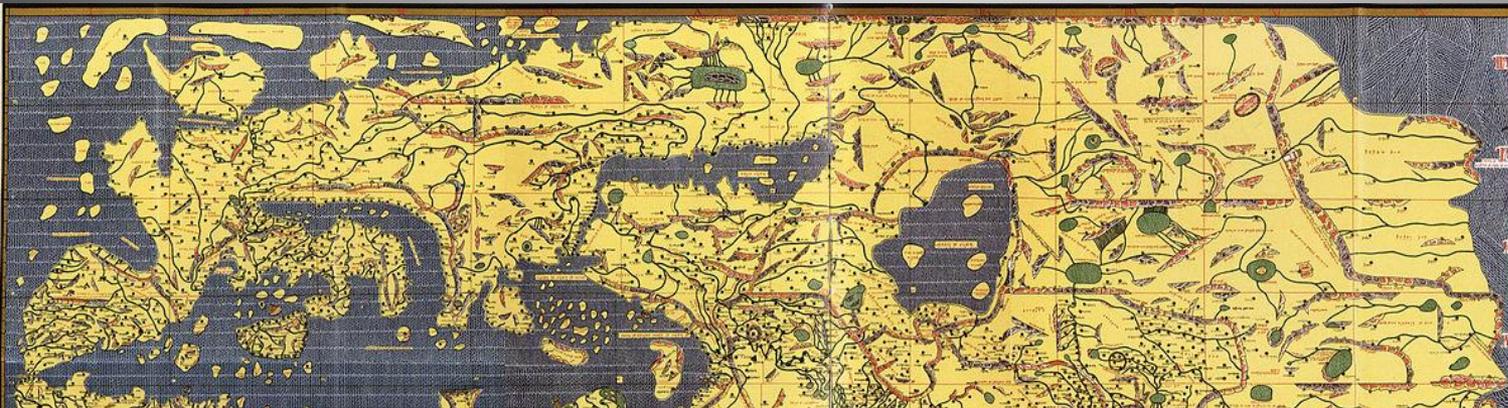
Once upon a  
time...



# Claudius Ptolemy (AD 90 – AD 168)

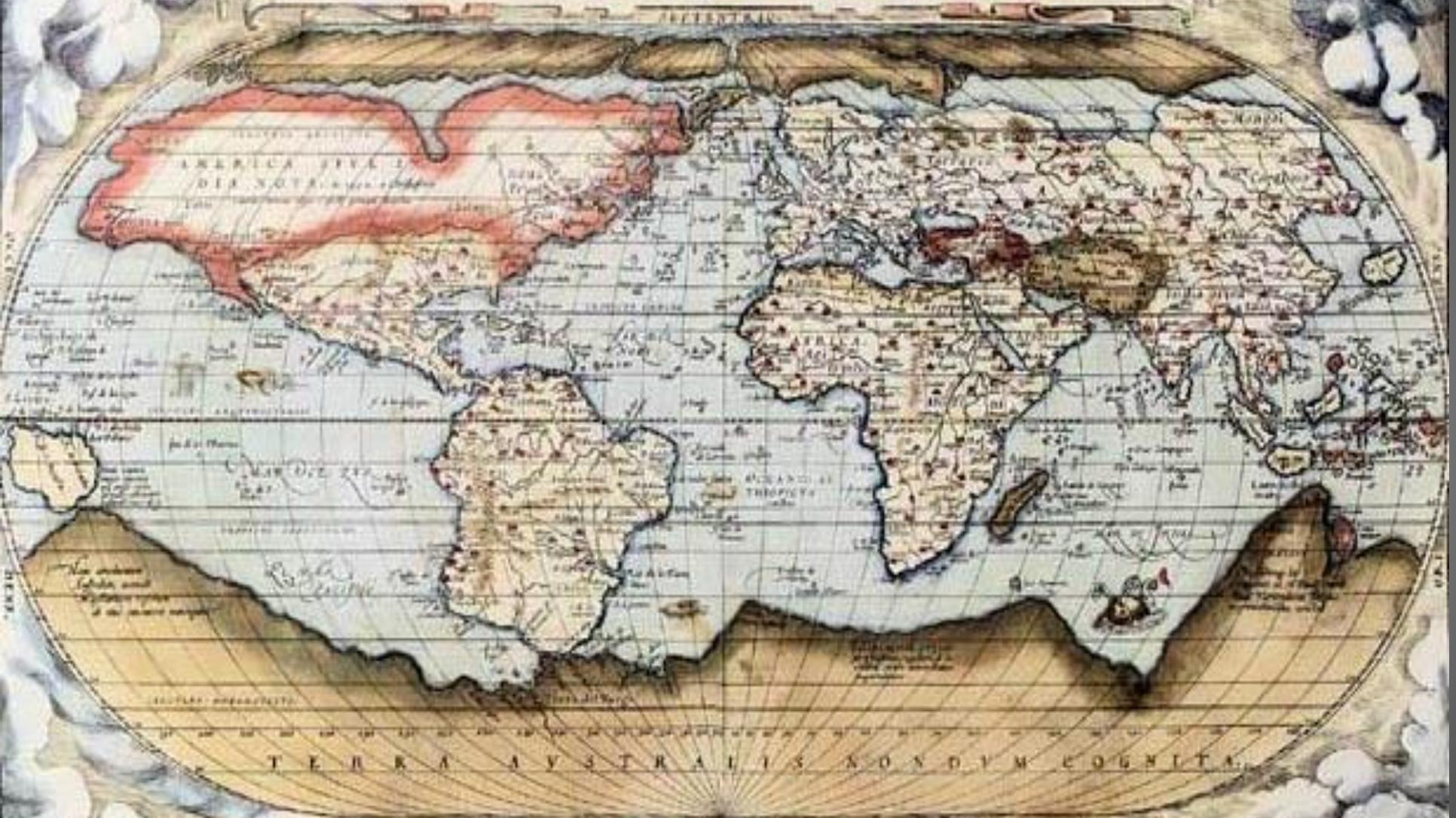


# Al Idrisi (1099–1165)





TYPVS ORBIS TERRARVM

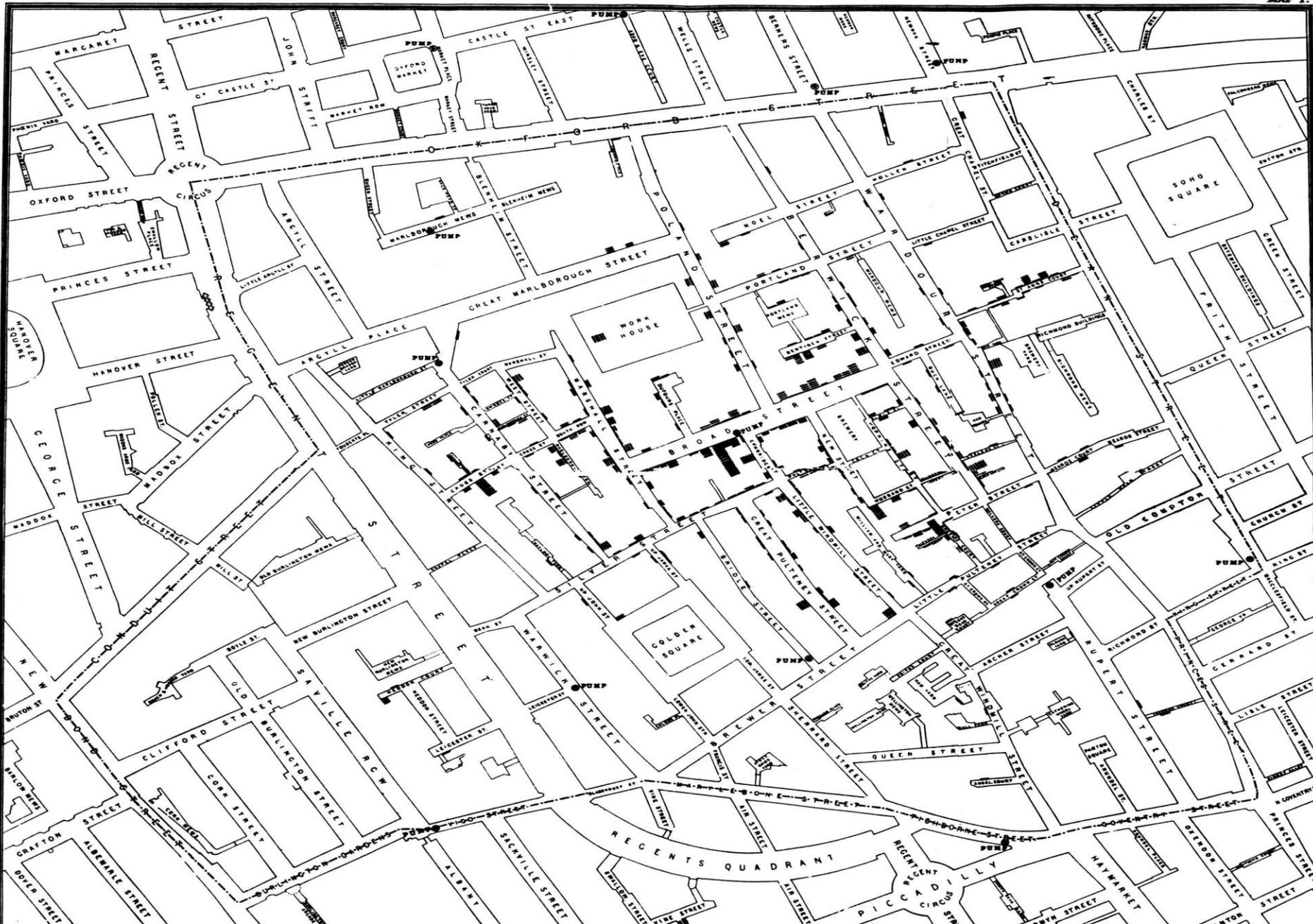


TERRA AVSTRALIS NONDVM COGNITA

QVID EI POTEST VIDERI MAGNUM IN REBVS HVMANIS, CVI AETERNITAS OMNIS, TOTIVSQUE MVNDI NOTA SIT MAGNITVDO. CICERO:



# Cholera cases in the London epidemic of 1854



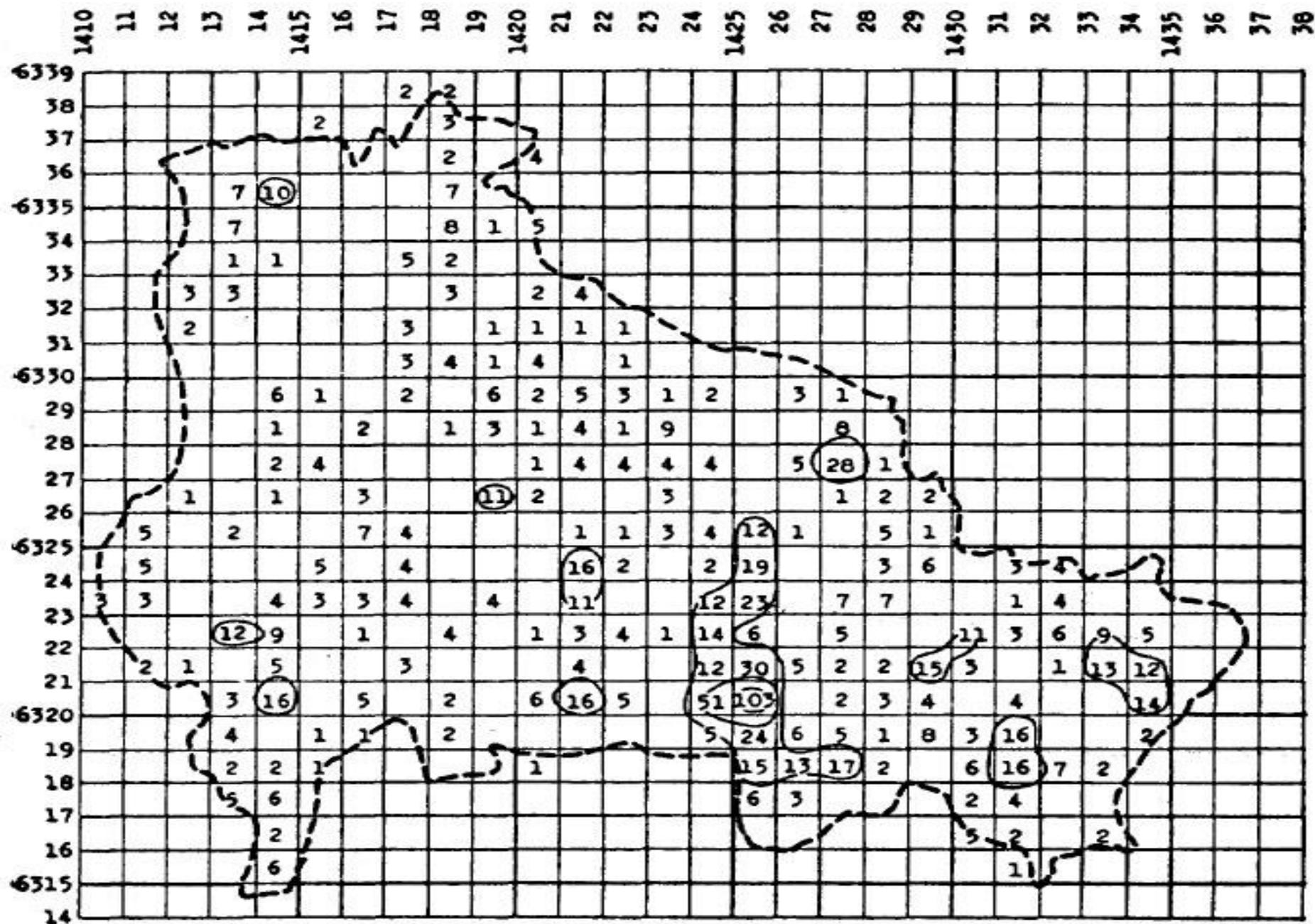


FIGURE 3—Children under 15 years of age in 1940.

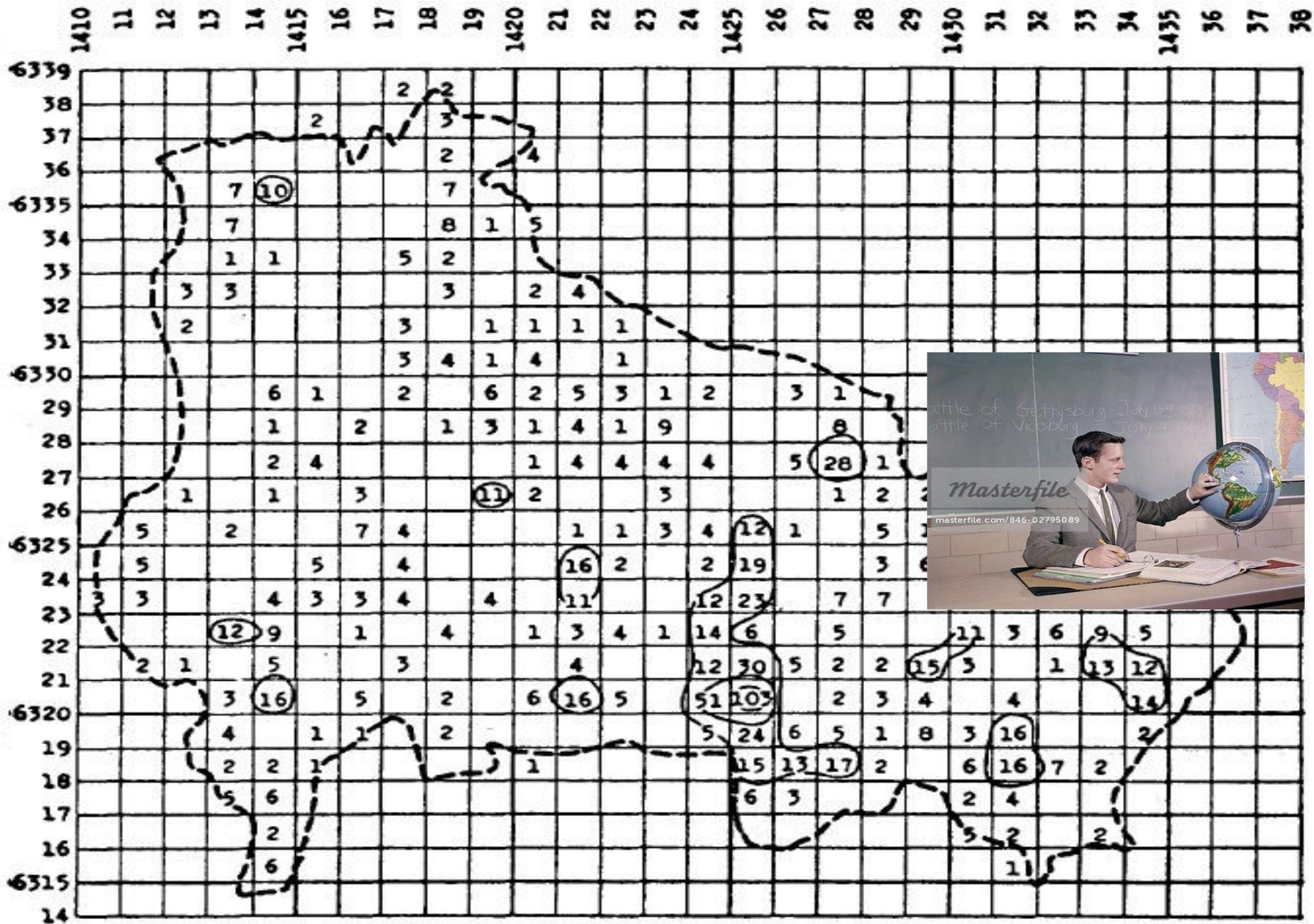
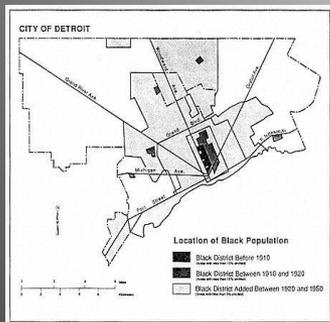
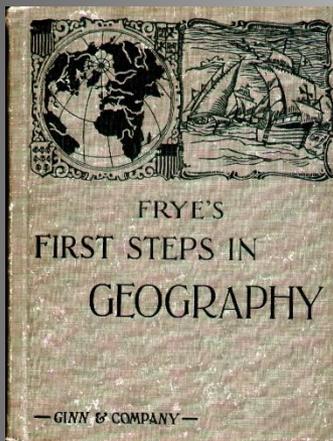
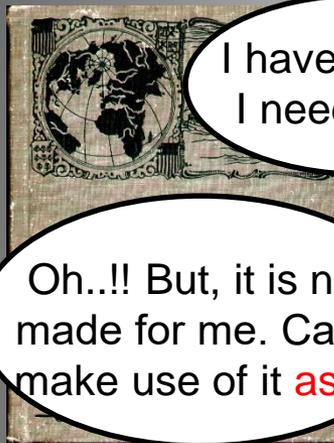


FIGURE 3—Children under 15 years of age in 1940.

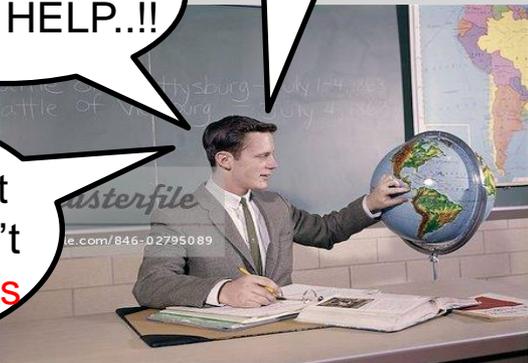




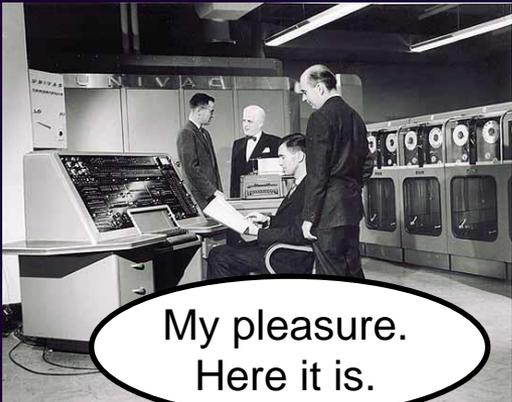
Cool **computer** technology..!!  
Can I use it in my application?



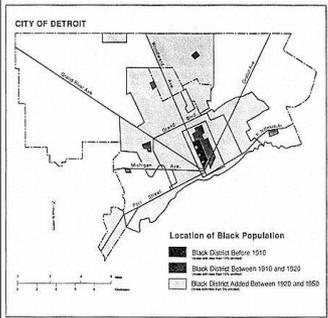
I have **BIG** data.  
I need **HELP**..!!

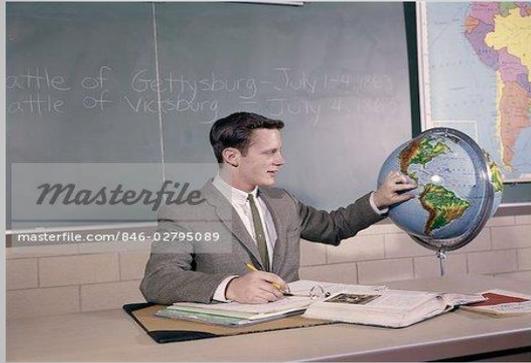


Oh..!! But, it is not made for me. Can't make use of it **as is**



My pleasure.  
Here it is.

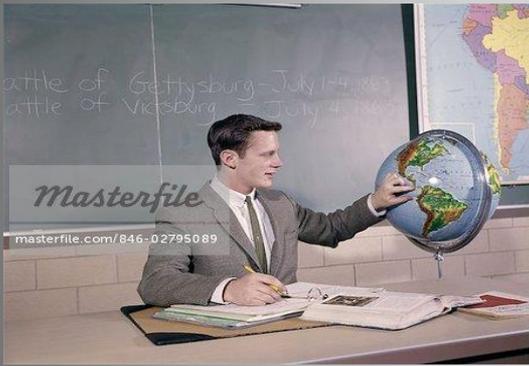




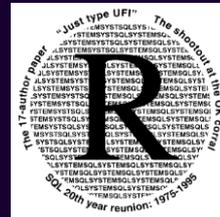
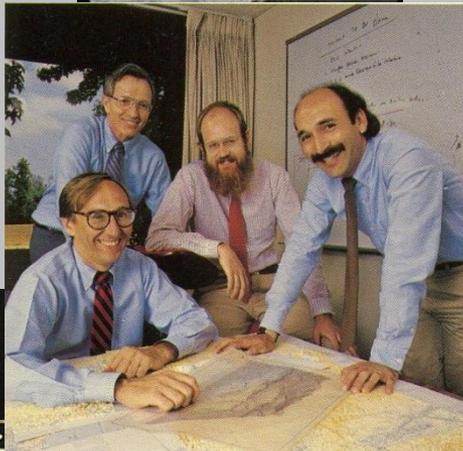
Kindly let me understand your needs

1969

Kindly let me get the technology you have



**ESRI**



# DATABASE MANAGEMENT SYSTEMS



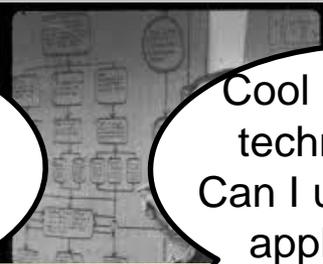
Informix

SQL

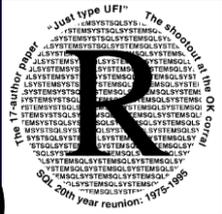




mmm...Let me check with my good friends there.



Cool **Database** technology..!! Can I use it in my application?



My pleasure. Here it is.

HELP..!! I have **BIG** data. Your technology is not helping me

Oh..!! But, it is not made for me. Can't make use of it **as is**



Informix

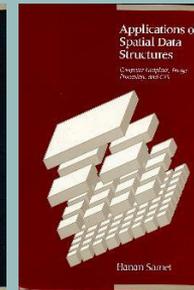
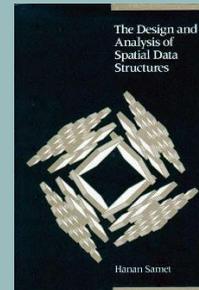
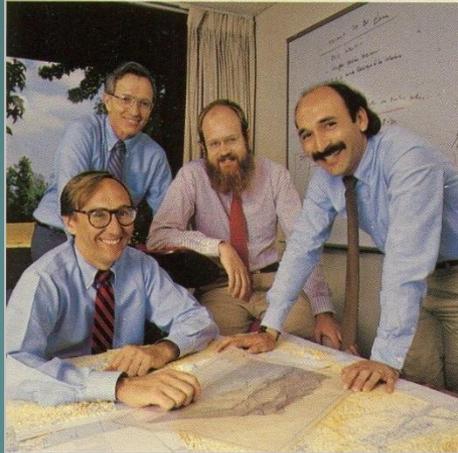
SQL

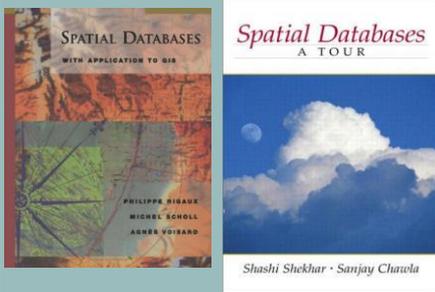
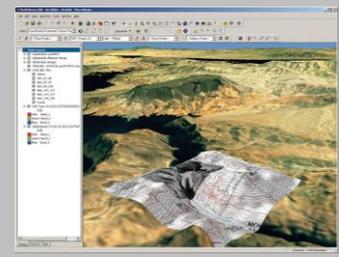
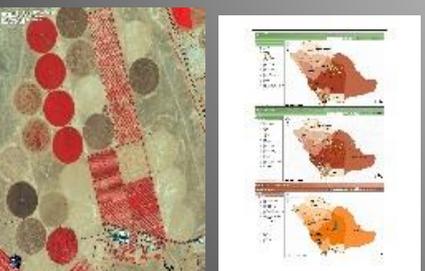
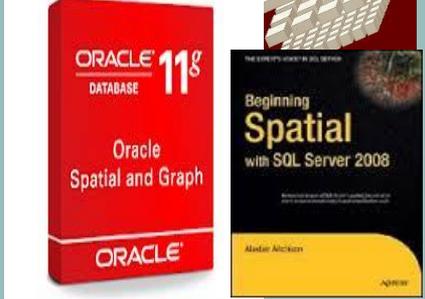
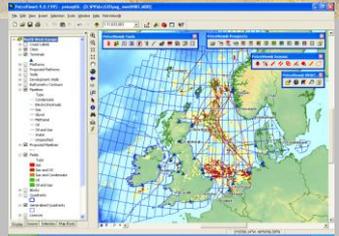
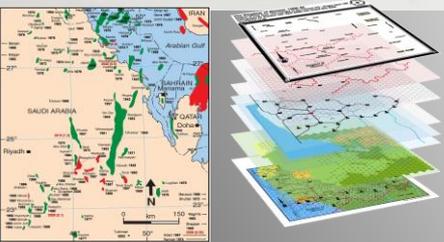
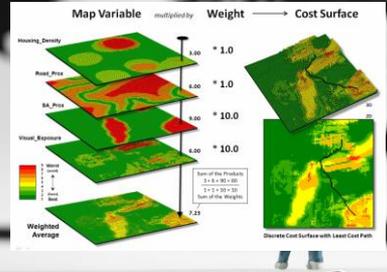


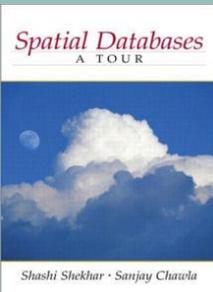
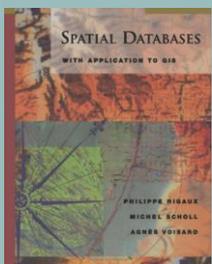
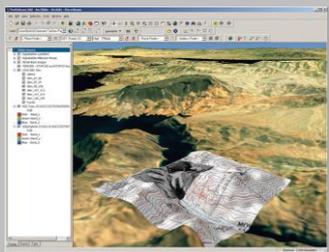
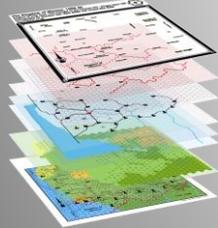
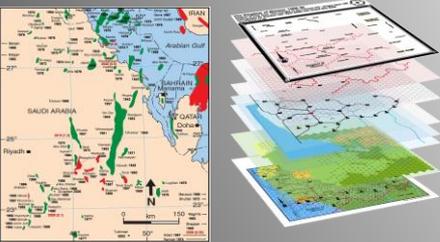
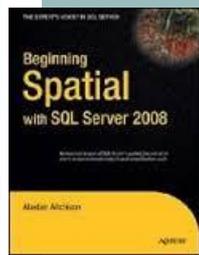
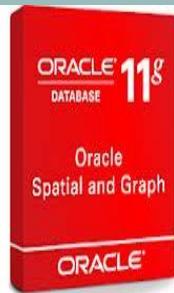
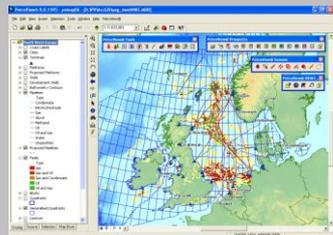
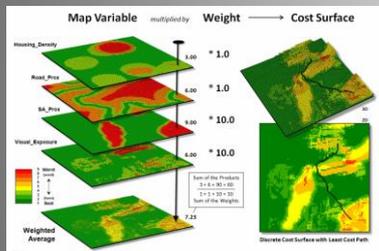


Kindly let me understand your needs

Kindly let me get the technology you have

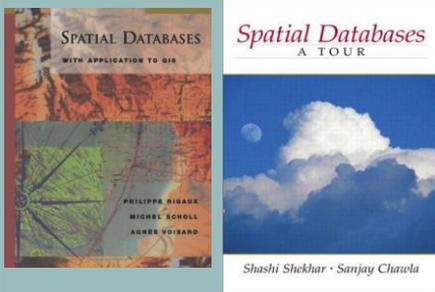
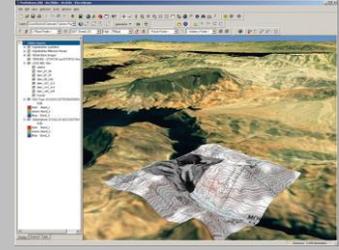
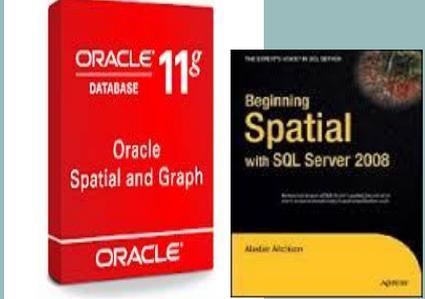
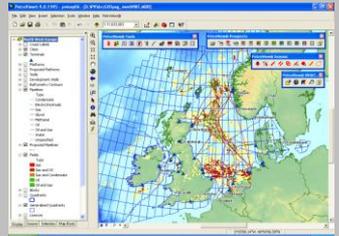








facebook **Map Reduce**



Let me check with my **other** good friends there.

Cool **Big Data** technology..!!  
Can I use it in my application?

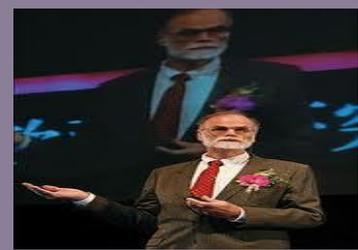
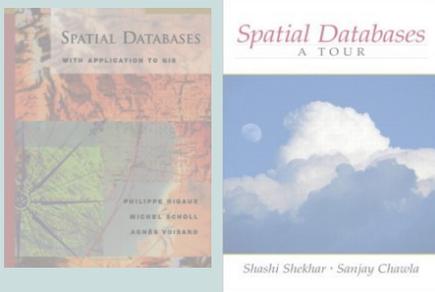
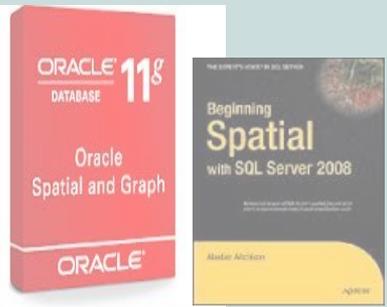
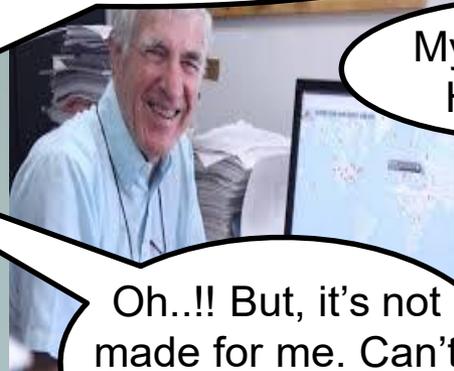
My pleasure.  
Here it is.

HELP..!! Again,  
I have **BIG** data.  
Your technology is  
not helping me

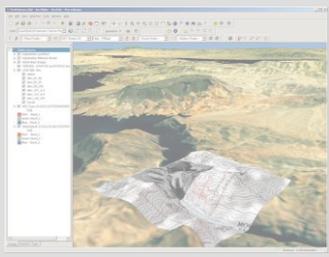
Sorry, seems like  
the DBMS  
technology cannot  
scale more

Oh..!! But, it's not  
made for me. Can't  
make use of it **as is**

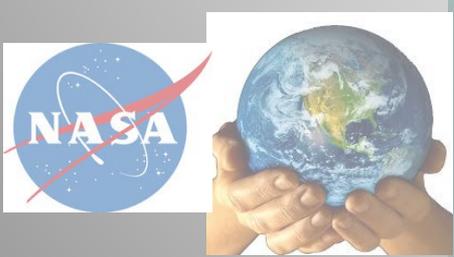
Google  
bing  
facebook *MapReduce*  
*hadoop*  
amazon web services™ *HIVE Spark*

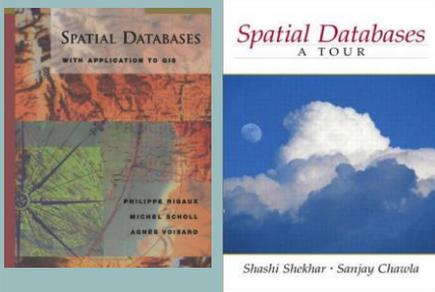
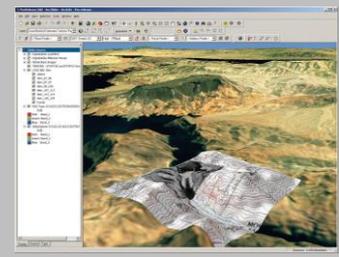
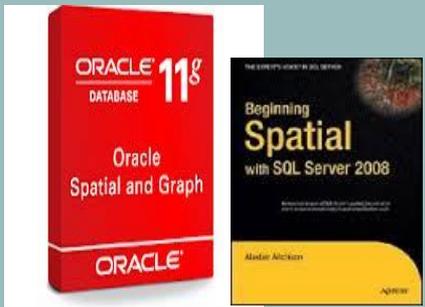
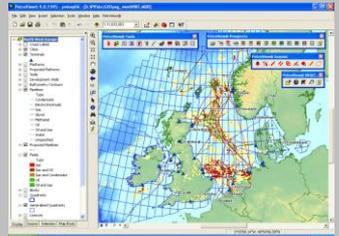


ORACLE®  
IBM DB2.  
Microsoft SQL Server™  
PostgreSQL  
MySQL  
SYBASE™ | An SAP Company



idrisi QGIS



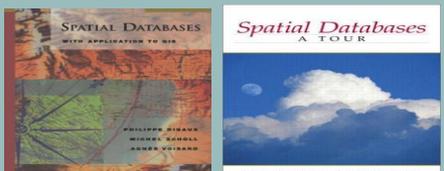
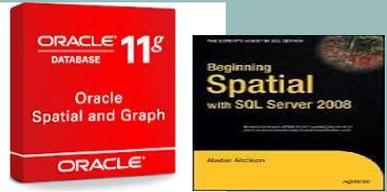
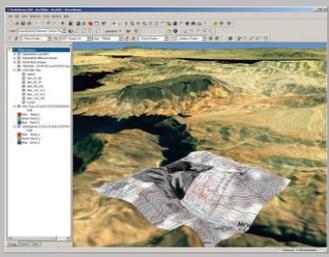
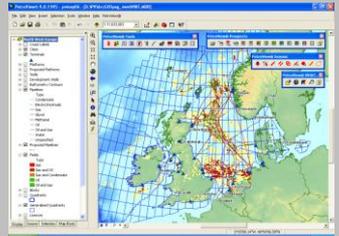




Kindly let me understand your needs



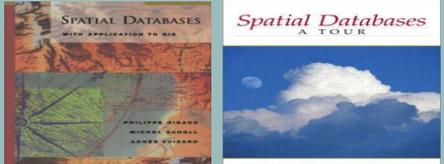
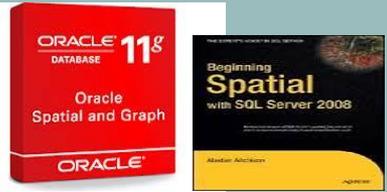
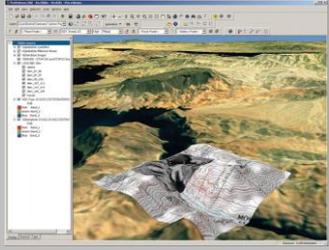
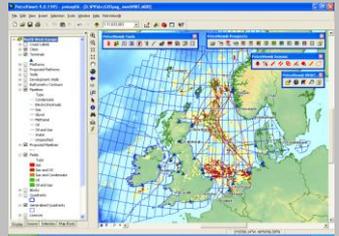
Kindly let me get the technology you have





# Big Spatial Data

Google  
bing  
twitter  
facebook *MapReduce*  
*hadoop*  
amazon web services™ *HIVE*  
*Spark*



# Tons of Spatial data out there...

twitter



Geotagged Microblogs



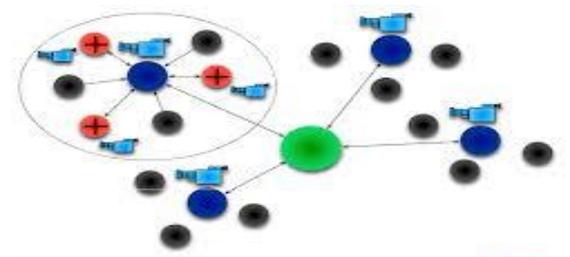
Geotagged Pictures



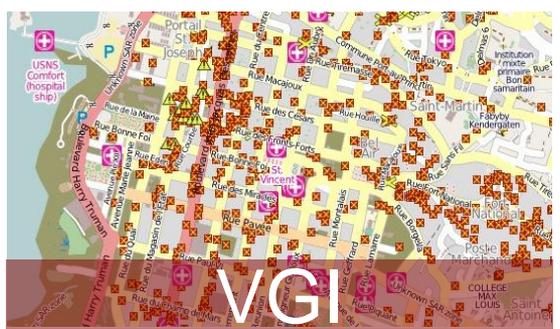
Medical Data



Smart Phones



Sensor Networks



VGI



Satellite Images



Traffic Data

# Spatial Data & Hadoop → SpatialHadoop



```
points = LOAD 'points' AS
  (id:int, x:int, y:int);
result = FILTER points BY
  x < xmax AND x >= xmin AND
  y < ymax AND y >= ymin;
```

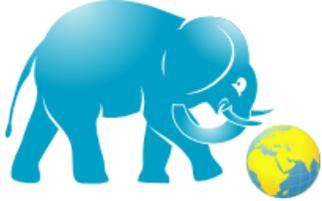
```
points = LOAD 'points' AS
  (id:int, location:point);
result = FILTER points BY
  Overlap(location, rectangle
  (xmin, ymin, xmax, ymax));
```



Takes 193 seconds



Finishes in 2 seconds



# Spatial Hadoop



KNN  
Point  
IsOverlap  
Rectangle  
DistanceTo

Spatial Language

Spatial Indexes

Spatial Operations

Visualization



80,000 downloads  
in one year



Conducted more than seven  
keynotes, tutorials, and invited talks

## Industry

## Academia



### Students Projects



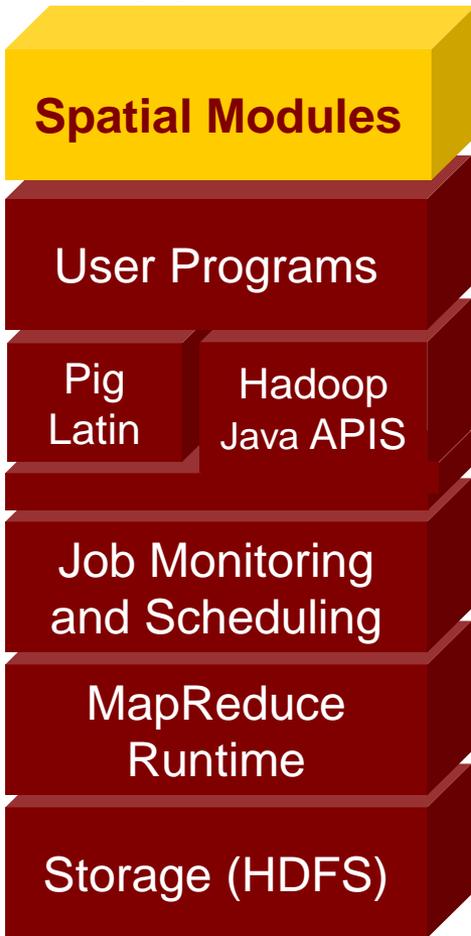
### Collaboration



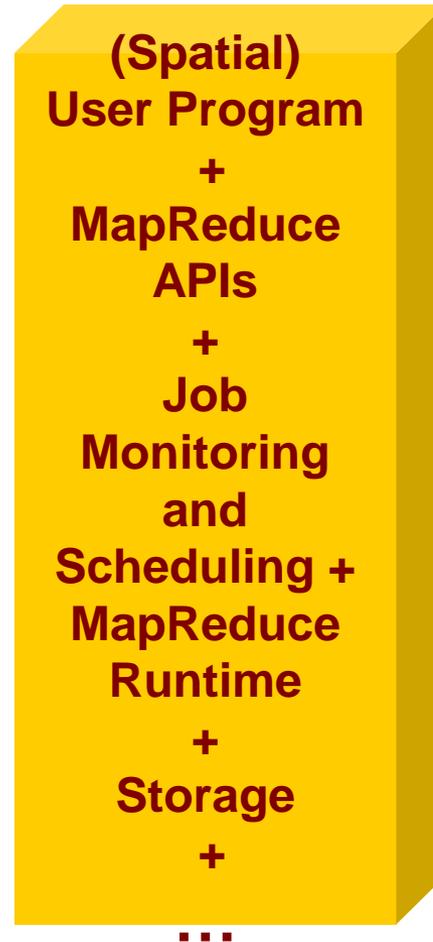
>500GB public datasets for  
benchmarking and testing

# The Built-in Approach of SpatialHadoop

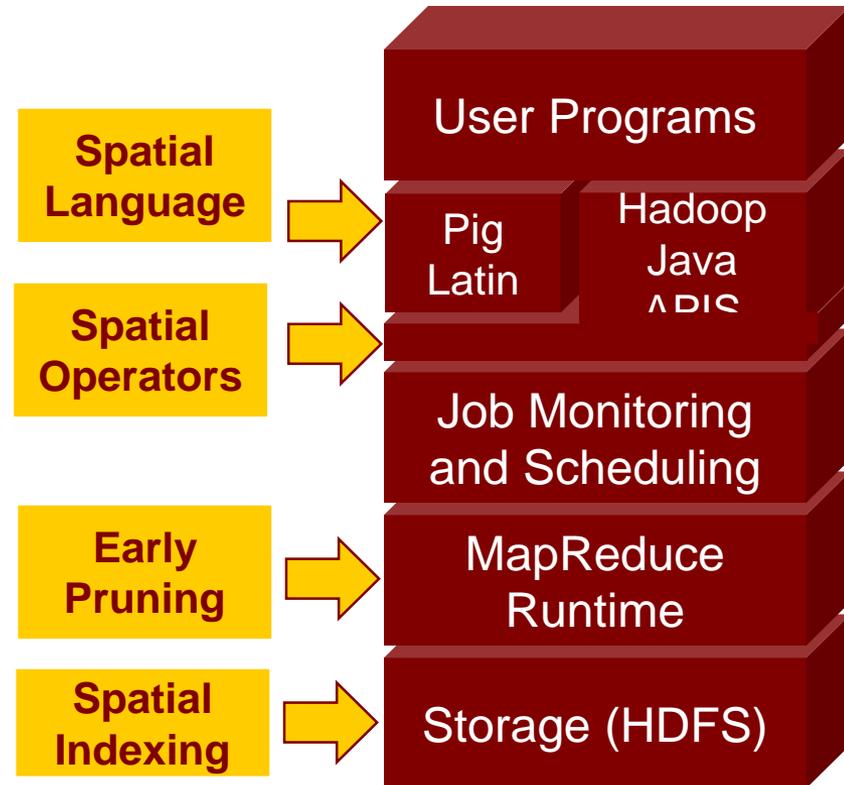
## The On-top Approach



## From Scratch Approach



## The Built-in Approach (SpatialHadoop)



# Agenda

- The ecosystem of SpatialHadoop
  - Motivation
  - Internal system design
  - Applications
  - Related work
  - Performance results
- Open Research Problems

# SpatialHadoop Architecture

**Applications:** SHAHED [ICDE'15] – MNTG [SSTD'13, ICDE'14]  
TAREEG[SIGMOD'14, SIGSPATIAL'14]



VLDB'13  
ICDE'15

**Language**  
Pigeon [ICDE'14]



**Visualization**  
[VLDB'15, ICDE'16]

**Operations**

Basic operations – CG\_Hadoop  
[SIGSPATIAL'13]

**MapReduce**

Spatial File Splitter  
Spatial Record Reader

**Indexing**

Grid – R-tree – R+-tree – Quad tree  
[VLDB'15]

ST-Hadoop

# Indexing

**Applications:** SHAHED [ICDE'15] – MNTG [SSTD'13, ICDE'14]  
TAREEG[SIGMOD'14, SIGSPATIAL'14]



VLDB'13  
ICDE'15

**Language**

Pigeon [ICDE'14]



**Visualization**

[VLDB'15, ICDE'16]

**Operations**

Basic operations – CG\_Hadoop  
[SIGSPATIAL'13]

**MapReduce**

Spatial File Splitter  
Spatial Record Reader

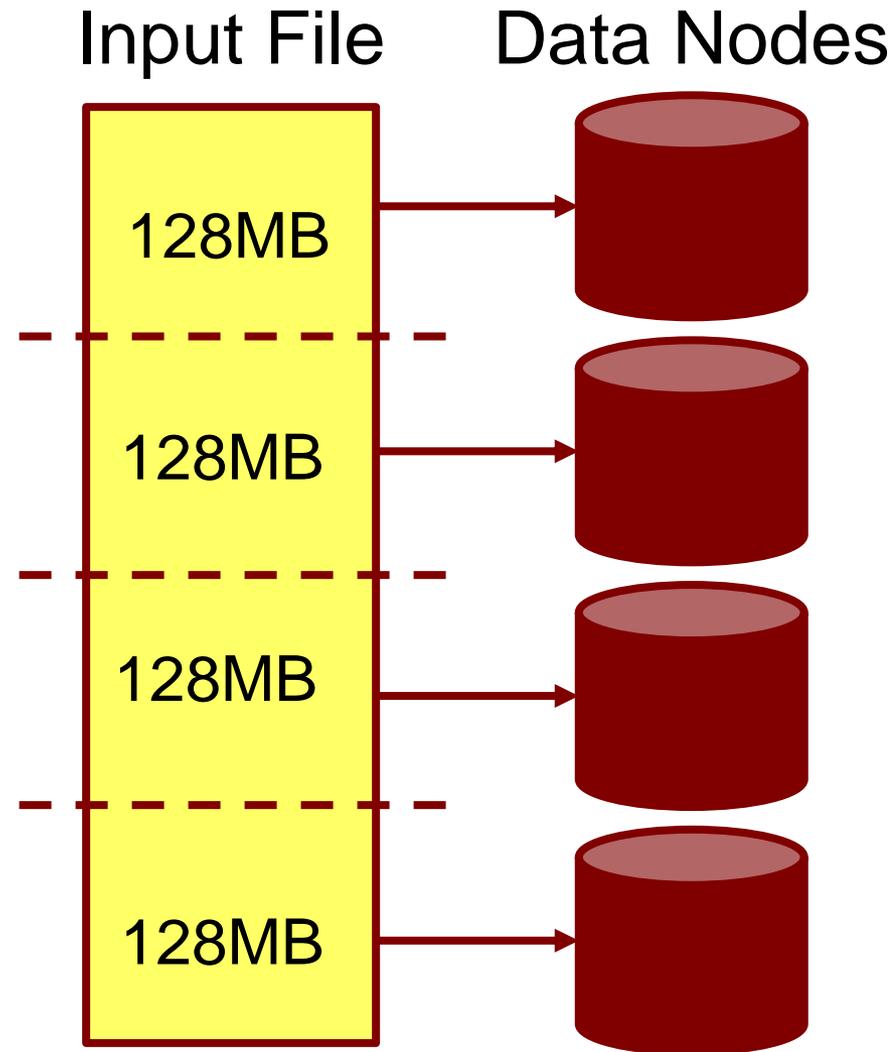
**Indexing**

Grid – R-tree – R+-tree – Quad tree  
[VLDB'15]

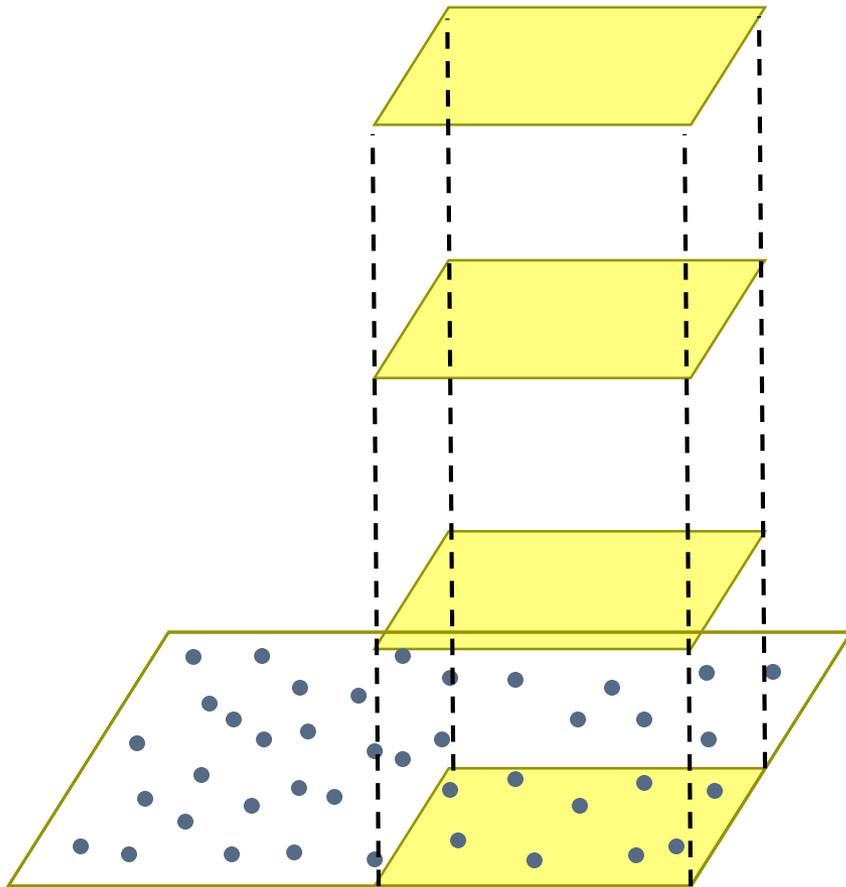
ST-Hadoop

# Data Loading in Hadoop

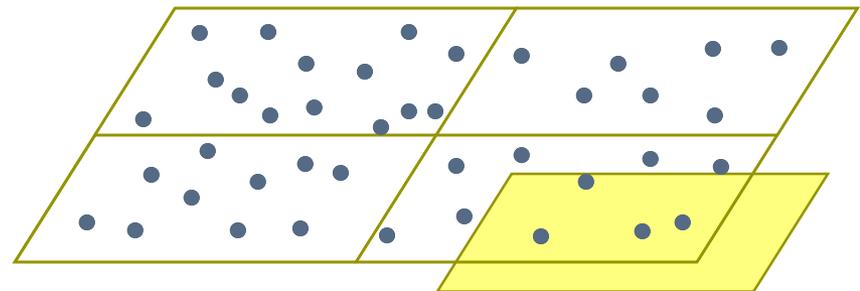
- Blindly chops down a big file into 128MB chunks
- Values of records are not considered
- Relevant records are typically assigned to two different blocks
- HDFS is too restrictive where files cannot be modified



# Spatial Distributed File System

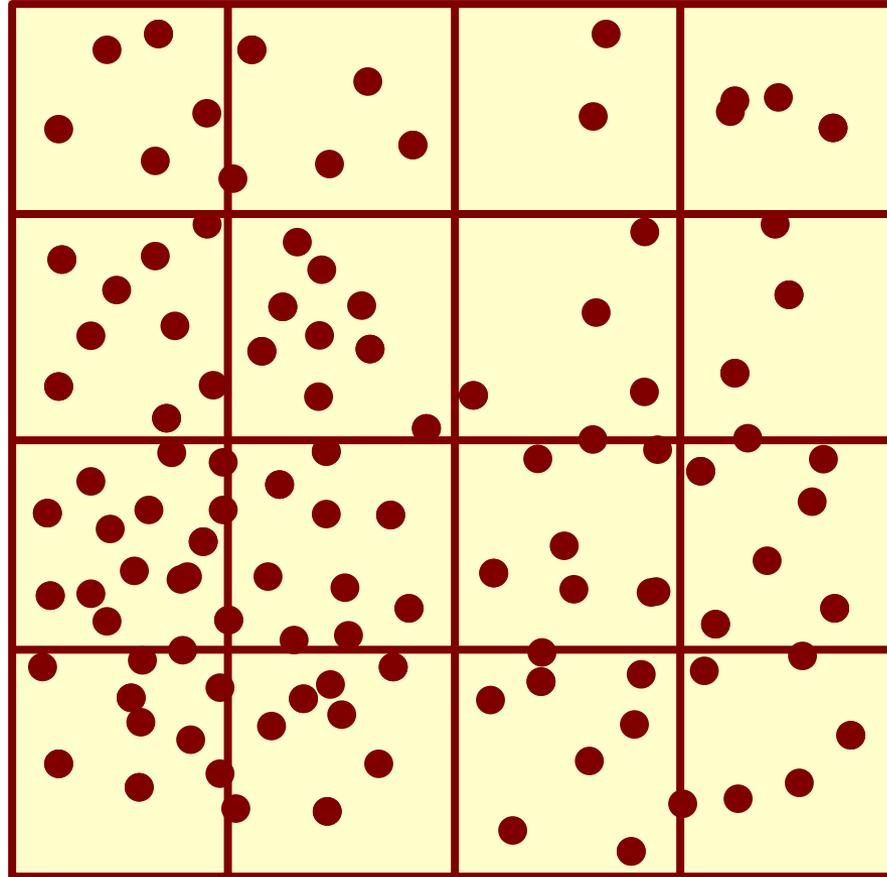


Default Partitioning



Spatial Partitioning

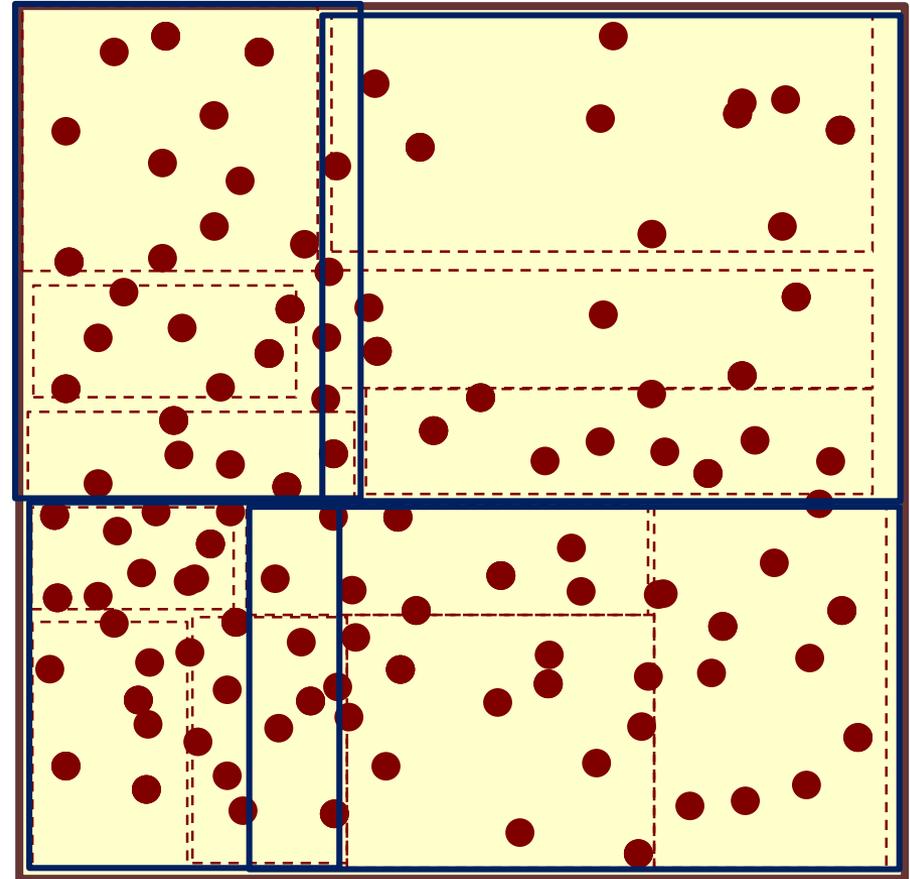
# Uniform Grid



Works only for uniformly distributed data

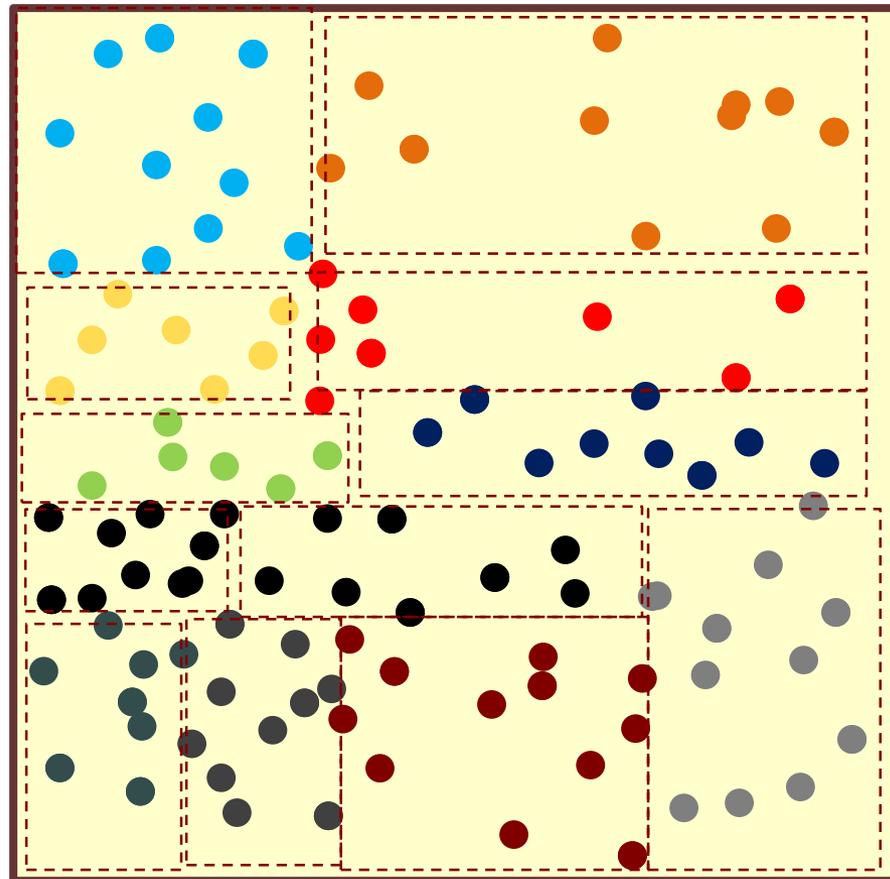
# R-tree

- › Read a sample
- › Bulk load the sample into an R-tree
  - › Leaf node capacity  $C$ 
$$C = \frac{k \cdot B}{|R|(1 + \alpha)}$$
    - ›  $k$ : Sample size
    - ›  $B$ : HDFS Block capacity
    - ›  $|R|$ : Input size
    - ›  $\alpha$ : Index overhead
- › Use MBR of leaf nodes as partition boundaries



# R-tree

- › Read a sample
- › Bulk load the sample into an R-tree
  - › Leaf node capacity  $C$ 
$$C = \frac{k \cdot B}{|R|(1 + \alpha)}$$
    - ›  $k$ : Sample size
    - ›  $B$ : HDFS Block capacity
    - ›  $|R|$ : Input size
    - ›  $\alpha$ : Index overhead
- › Use MBR of leaf nodes as partition boundaries
- › Partition the data



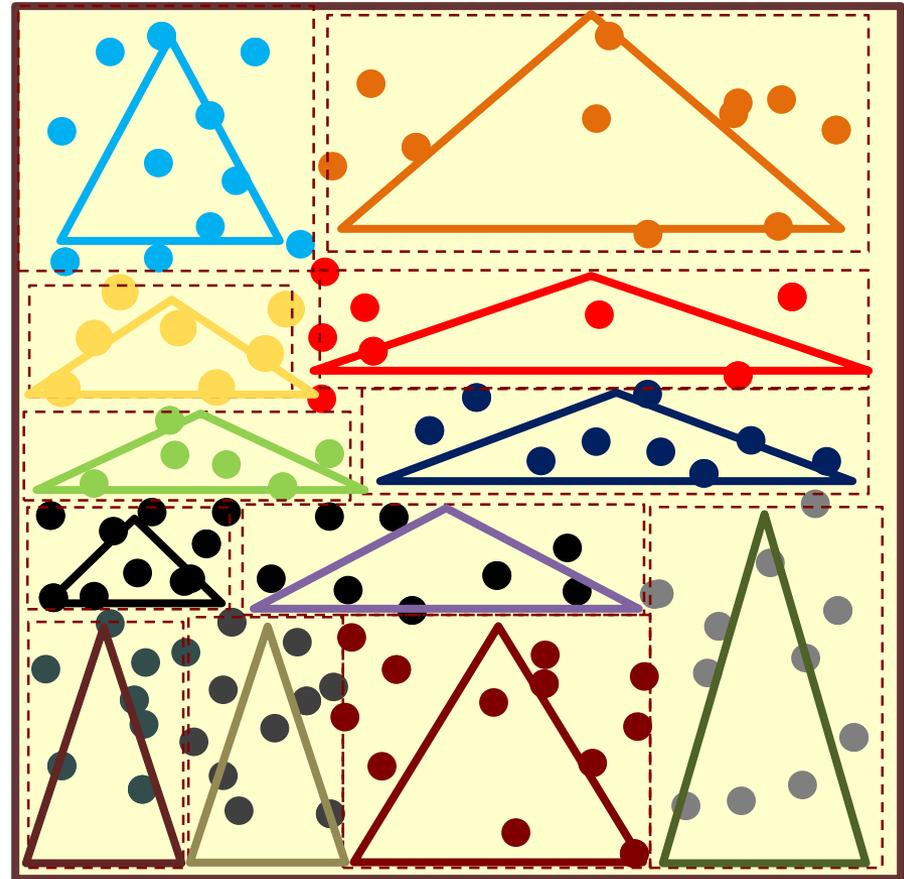
# R-tree

- › Read a sample
- › Bulk load the sample into an R-tree

- › Leaf node capacity  $C$

$$C = \frac{k \cdot B}{|R|(1 + \alpha)}$$

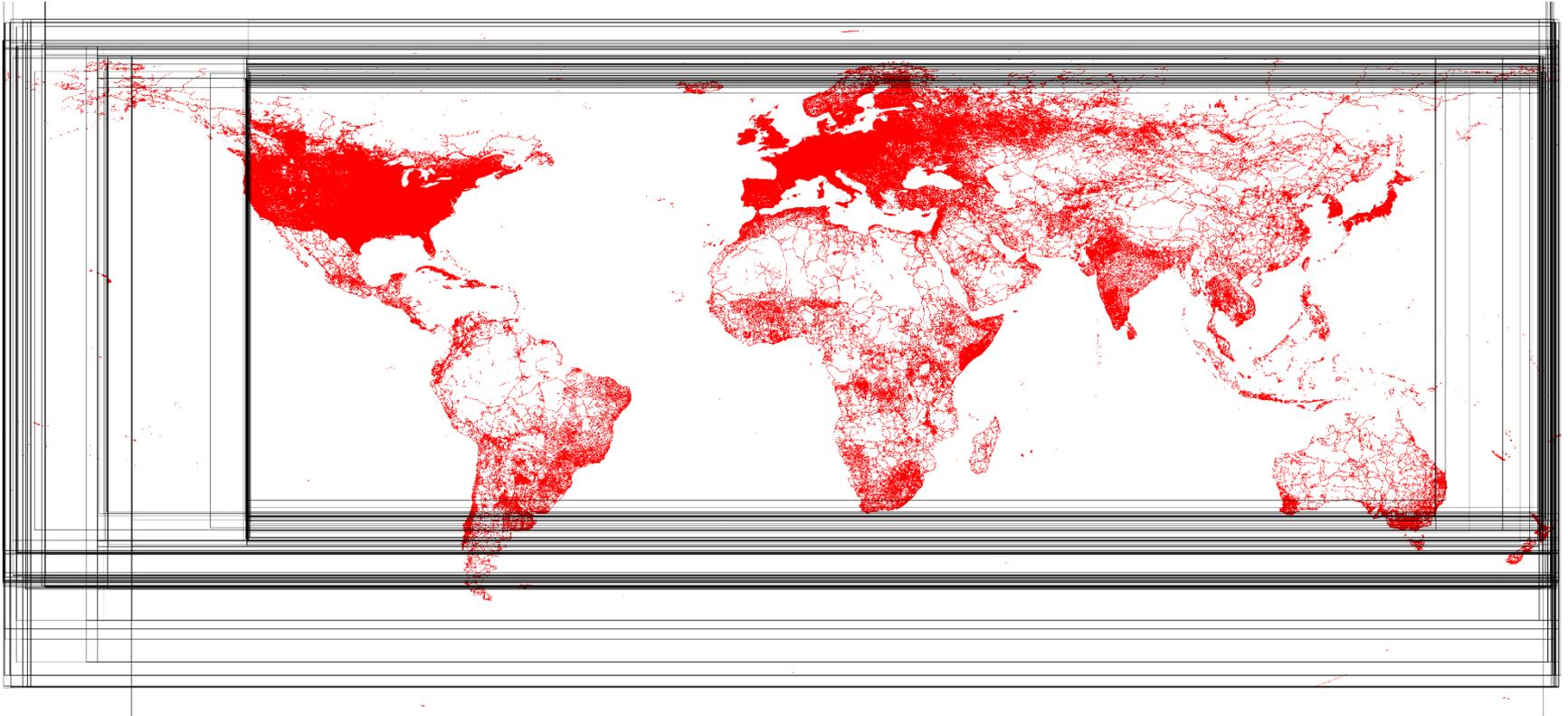
- ›  $k$ : Sample size
- ›  $B$ : HDFS Block capacity
- ›  $|R|$ : Input size
- ›  $\alpha$ : Index overhead
- › Use MBR of leaf nodes as partition boundaries
- › Partition the data
- › Optional: Build R-tree Local indexes



# R-tree-based Index of a 400 GB road network



# Non-indexed Heap File



# Operations

**Applications:** SHAHED [ICDE'15] – MNTG [SSTD'13, ICDE'14]  
TAREEG[SIGMOD'14, SIGSPATIAL'14]



VLDB'13  
ICDE'15

**Language**

Pigeon [ICDE'14]



**Visualization**

[VLDB'15, ICDE'16]

**Operations**

Basic operations – CG\_Hadoop  
[SIGSPATIAL'13]

**MapReduce**

Spatial File Splitter  
Spatial Record Reader

**Indexing**

Grid – R-tree – R+-tree – Quad tree  
[VLDB'15]

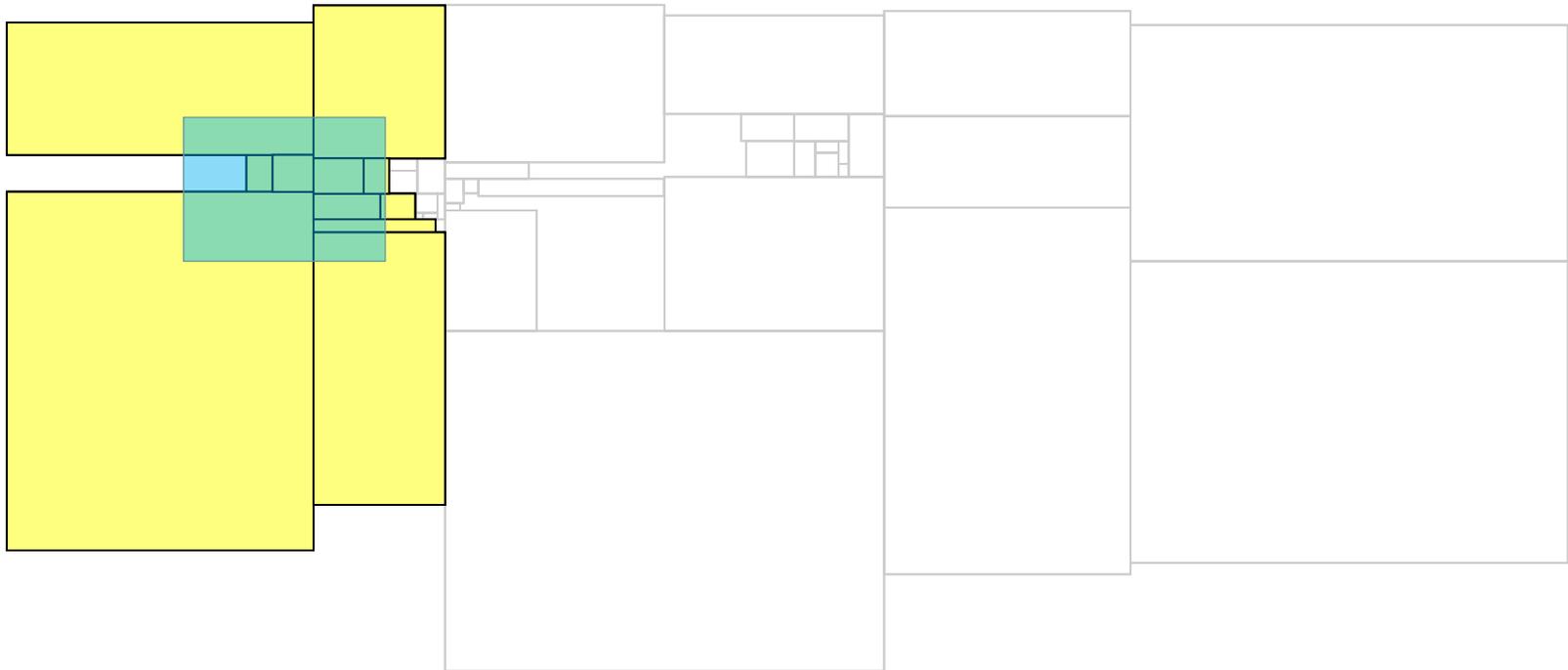
ST-Hadoop

# Operations Layer



- Basic Operations: e.g, Range query and KNN
- Spatial Join Operations
- Computational geometry operations: e.g., Polygon Union, Voronoi diagram, Delaunay Triangulation, and Convex Hull
- User-defined operations: e.g., kNN join

# Range Query



Use **local indexes** to find matching records

Use the **global index** to prune disjoint partitions

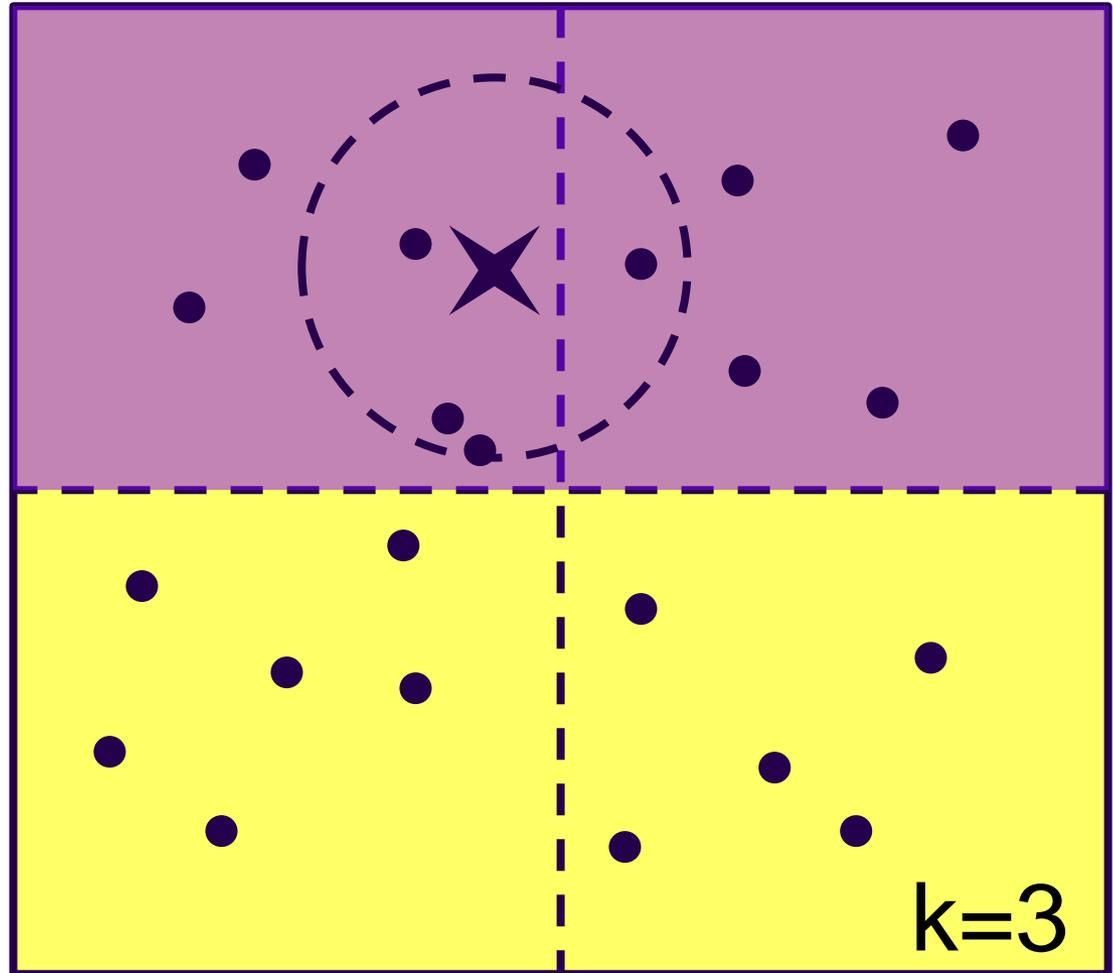
# KNN over Indexed Data

First iteration runs as before and result is tested for correctness

✗ Answer is incorrect

Second iteration processes other blocks that might contain an answer

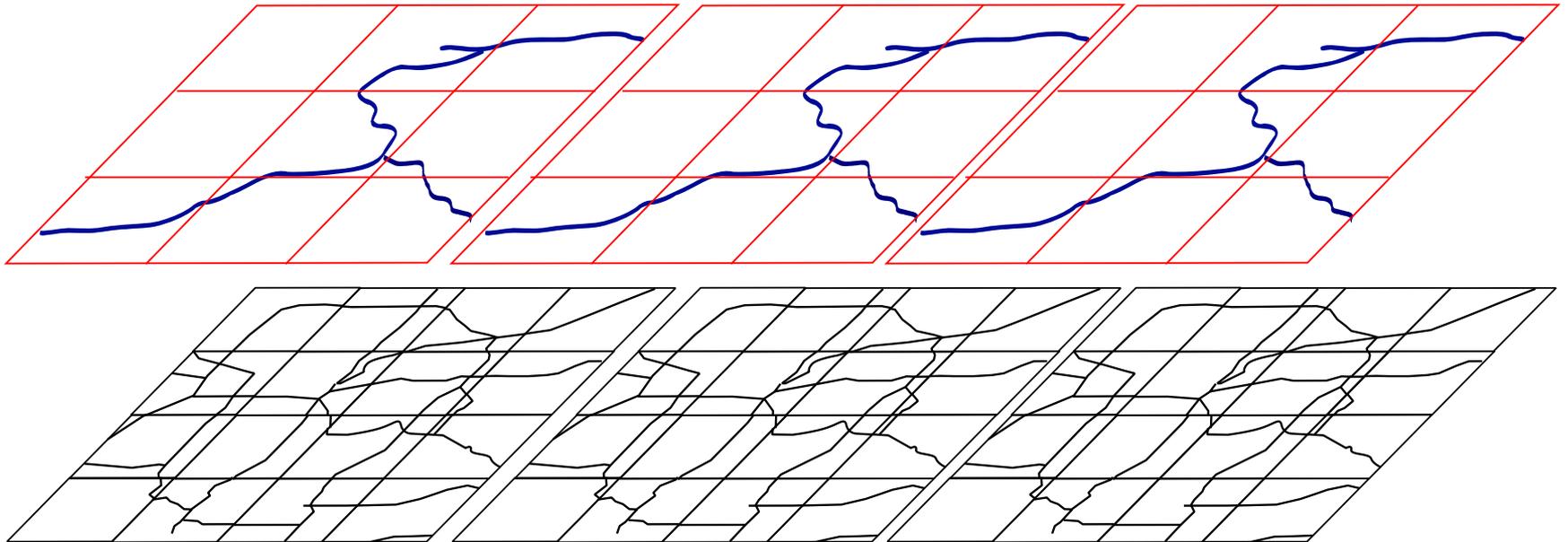
✓ Answer is correct



# Spatial Join

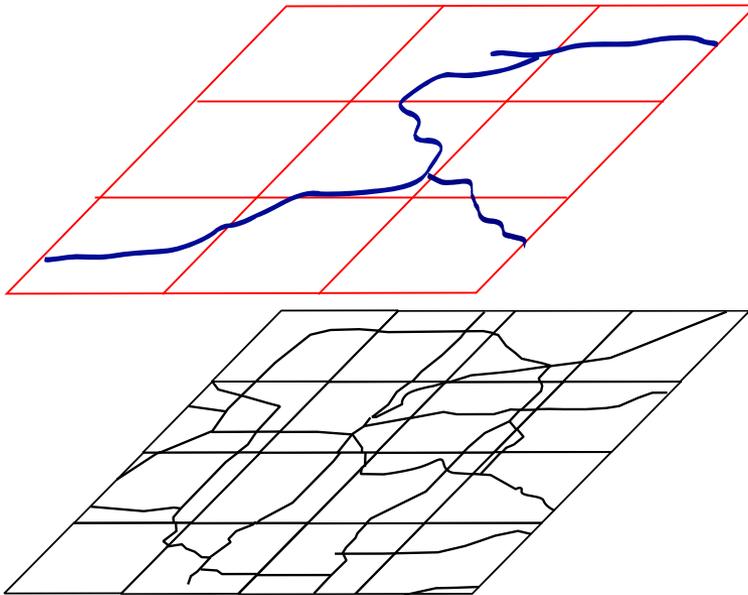
Join Directly

Partition – Join



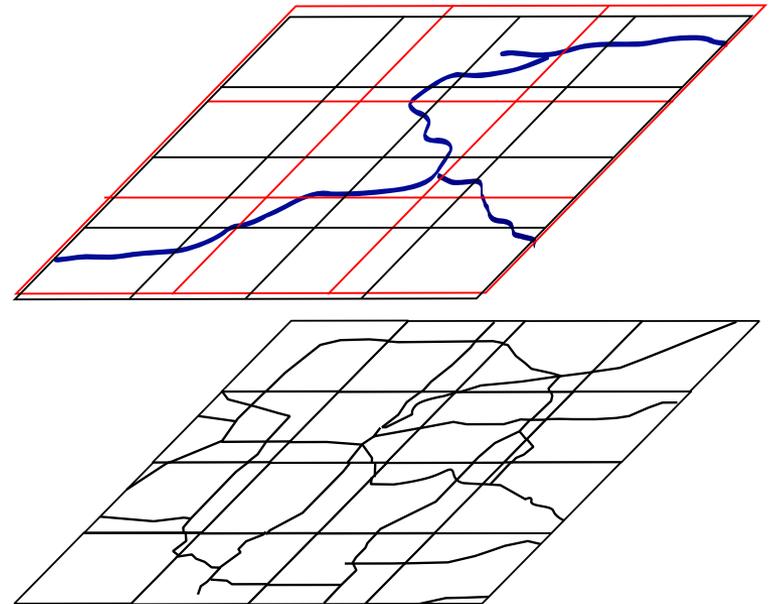
# Spatial Join

## Join Directly



Total of 36 overlapping pairs

## Partition – Join

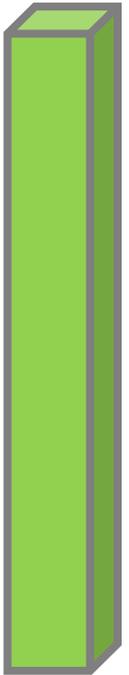


Only 16 overlapping pairs

# CG\_Hadoop

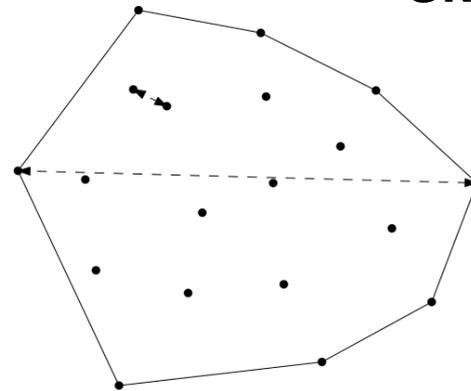
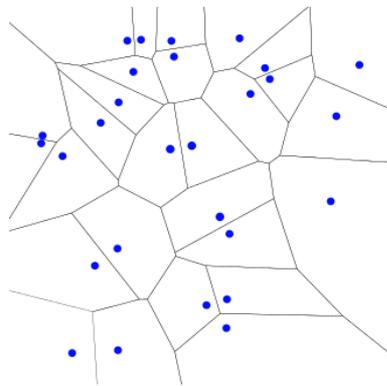
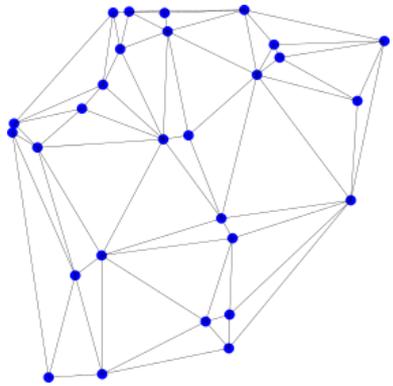


260x



**Polygon Union**

**Skyline**



**1x**  
Single Machine

**29x**  
Hadoop

Spatial Hadoop

**Delaunay Triangulation**

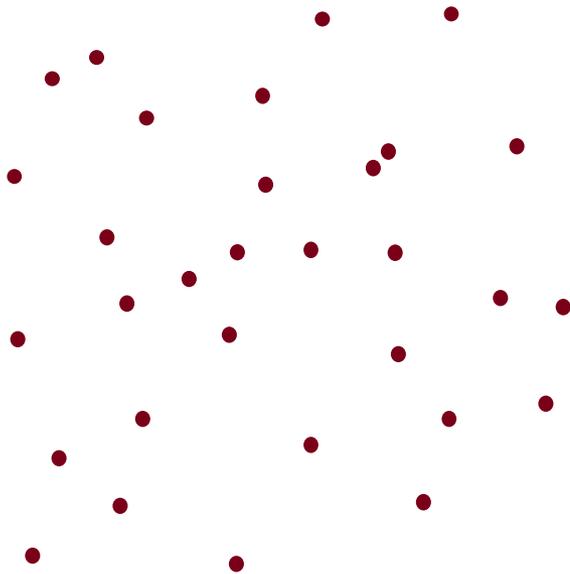
**Voronoi Diagram**

**Convex Hull  
Farthest/closest pair**

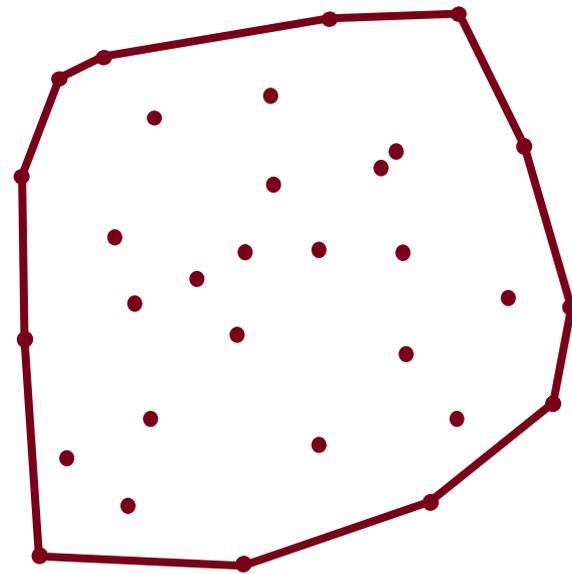
# Convex Hull

Find the minimal convex polygon that contains all points

Input



Output

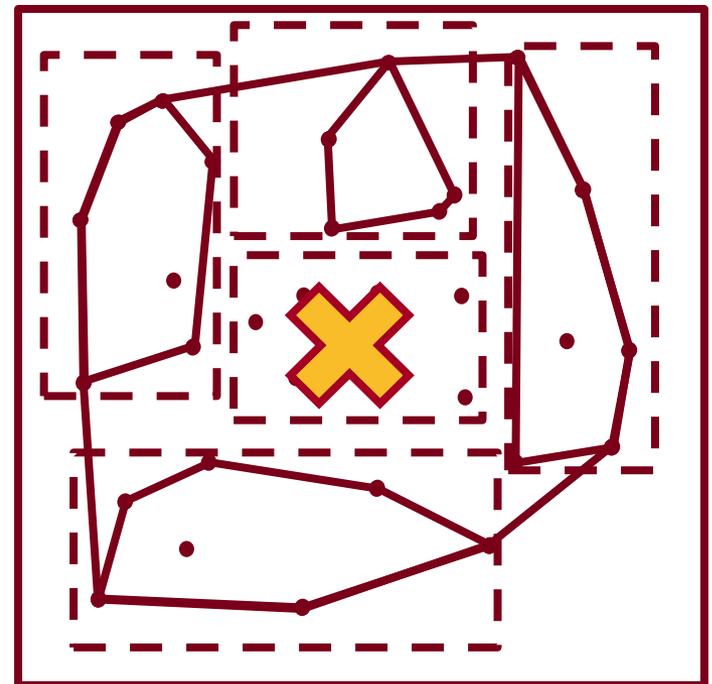
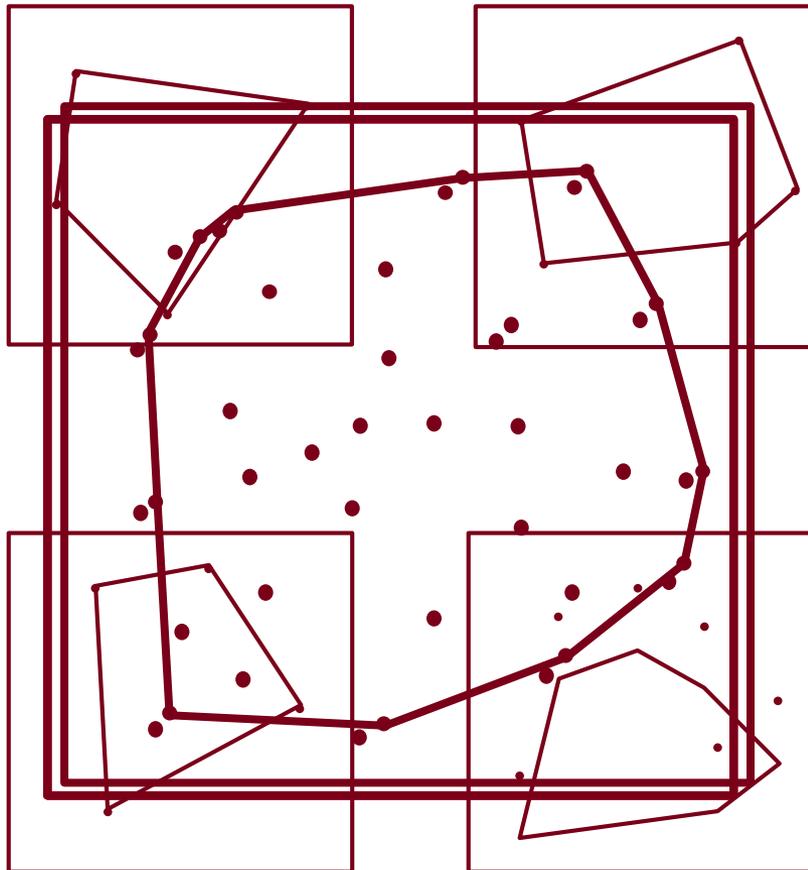


# Convex Hull in CG\_Hadoop

Hadoop

SpatialHadoop

- ① Partition
- ② Pruning
- ③ Local hull
- ④ Global hull



# Advanced Analytics

(Ongoing work)

Partitioning

Local VD

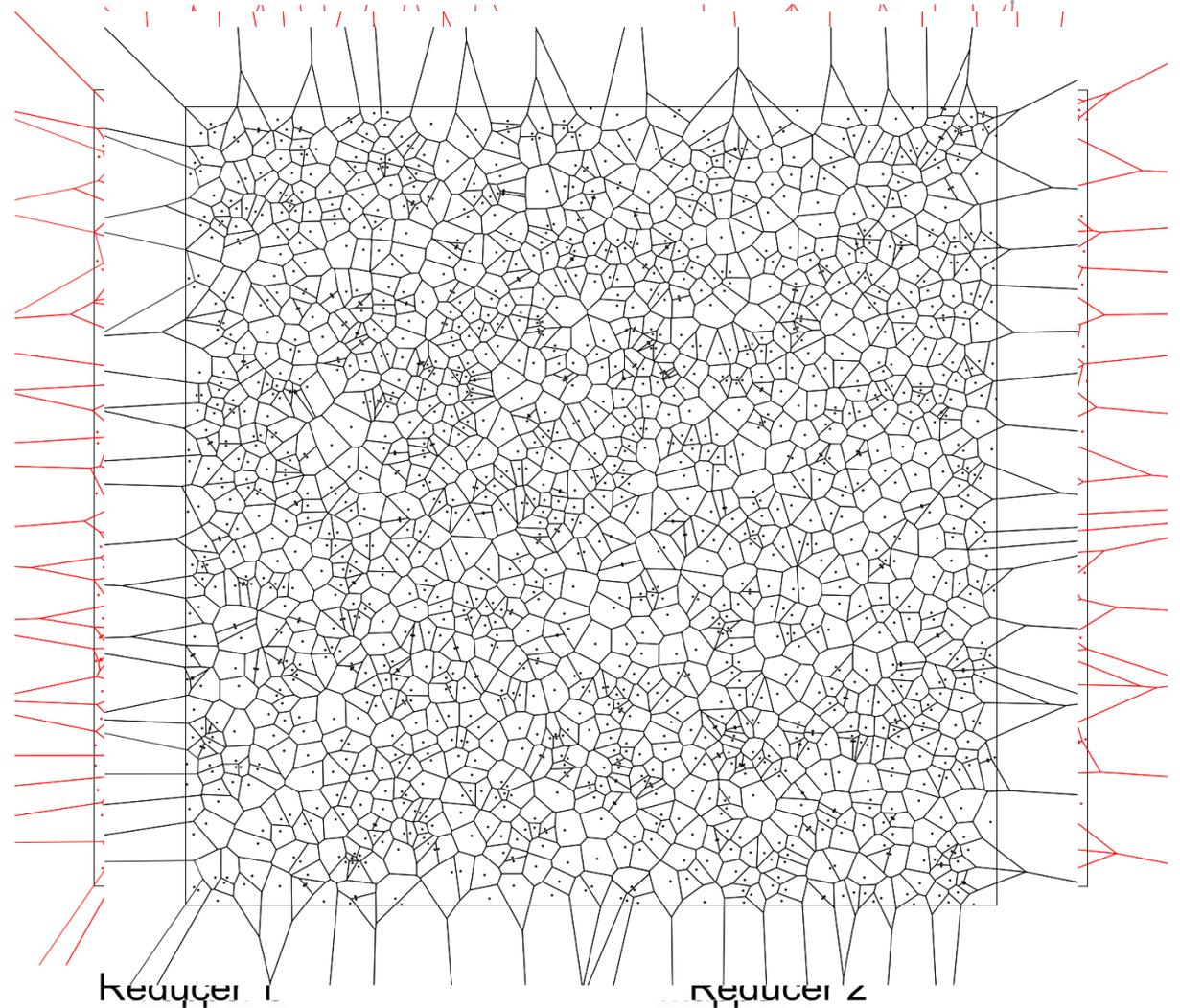
Pruning

Vertical Merge

Pruning

Horizontal Merge

Final output



# Applications

**Applications:** SHAHED [ICDE'15] – MNTG [SSTD'13, ICDE'14]  
TAREEG[SIGMOD'14, SIGSPATIAL'14]



VLDB'13  
ICDE'15

**Language**

Pigeon [ICDE'14]



**Visualization**

[VLDB'15, ICDE'16]

**Operations**

Basic operations – CG\_Hadoop  
[SIGSPATIAL'13,]

**MapReduce**

Spatial File Splitter  
Spatial Record Reader

**Indexing**

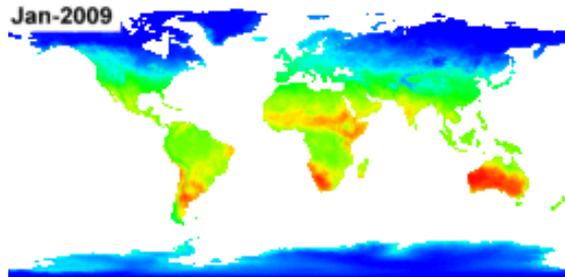
Grid – R-tree – R+-tree – Quad tree  
[VLDB'15]

ST-Hadoop

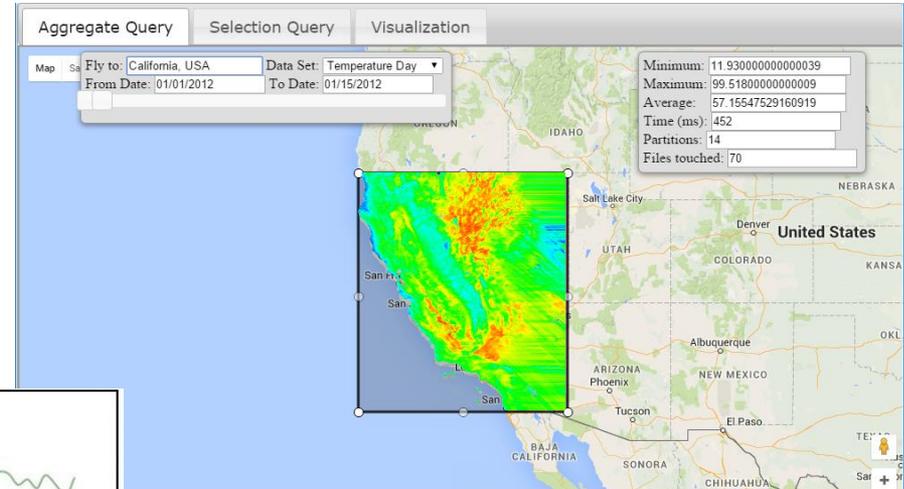
# SHAHERD – A system for querying and visualizing spatio-temporal satellite data



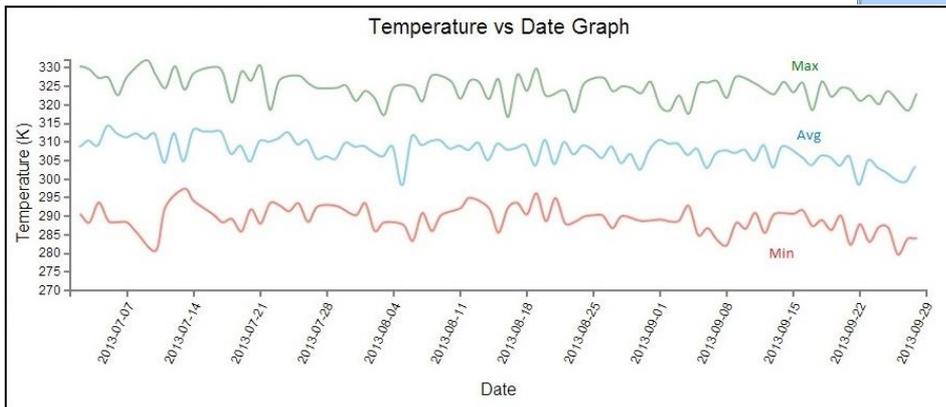
<http://shahed.cs.umn.edu/>



Visualize animated heat maps or still images



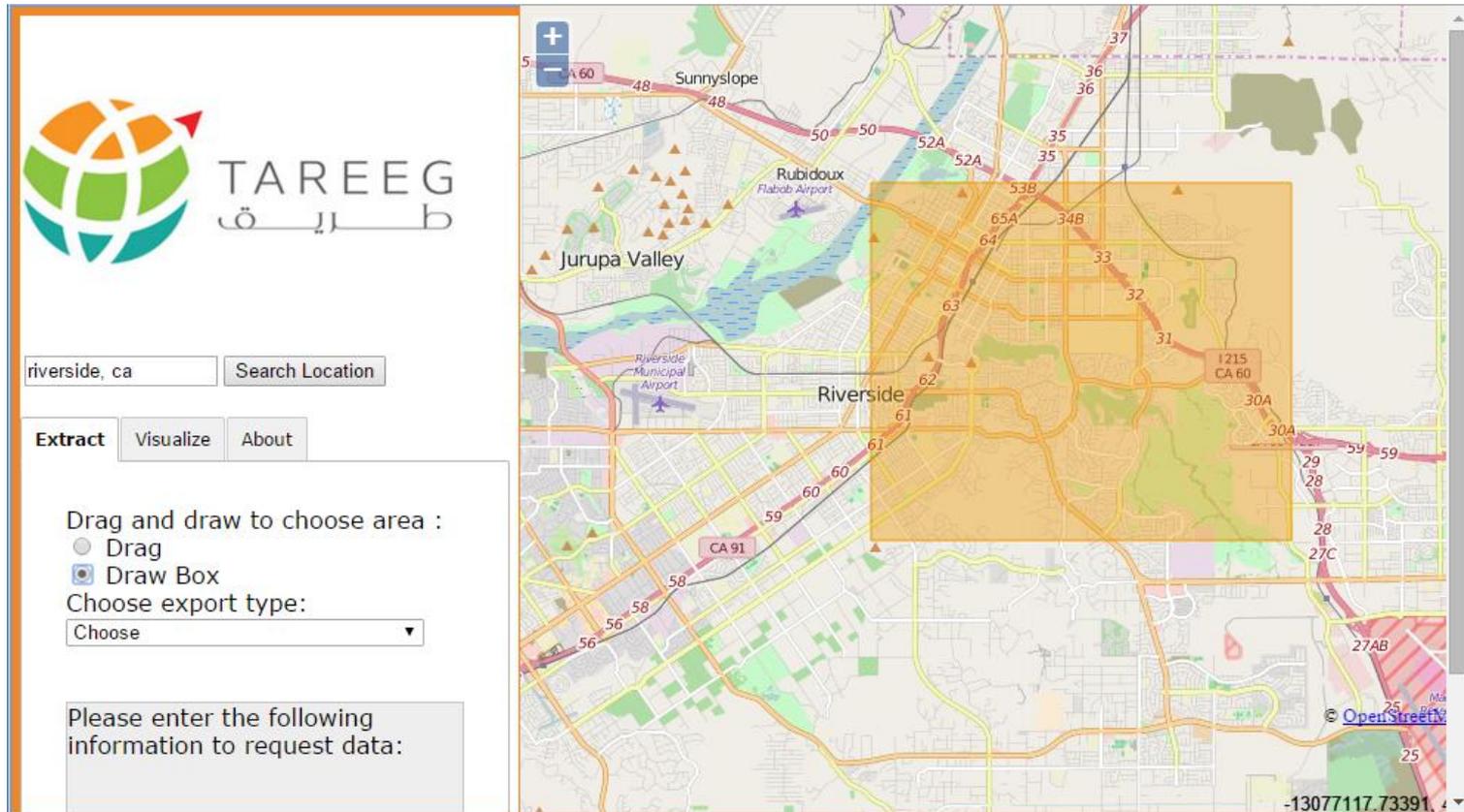
Run spatio-temporal selection and aggregate queries



A. Eldawy *et al.* "SHAHERD: A MapReduce-based System for Querying and Visualizing Spatio-temporal Satellite Data", **IEEE ICDE'15 (Best poster runner-up)**  
A. Eldawy *et al.* "A Demonstration of SHAHERD: A MapReduce-based System for Querying and Visualizing Satellite Data", **IEEE ICDE'15**

# TAREEG – Web-based extractor for OpenStreetMap data using MapReduce

<http://tareeg.net/>



The screenshot displays the TAREEG web interface. On the left, there is a search bar containing "riverside, ca" and a "Search Location" button. Below the search bar are three tabs: "Extract" (selected), "Visualize", and "About". Under the "Extract" tab, there are instructions: "Drag and draw to choose area :", followed by two radio buttons: "Drag" (unselected) and "Draw Box" (selected). Below this is a "Choose export type:" label and a dropdown menu currently set to "Choose". At the bottom left, there is a text box with the prompt "Please enter the following information to request data:". The main part of the interface is a map of Riverside, CA, showing a large orange rectangular selection box over the city center. The map includes labels for "Sunnyslope", "Rubidoux Flabob Airport", "Jurupa Valley", and "Riverside Municipal Airport". The map also shows various street names and highway markers like "CA 91" and "CA 60". The bottom right corner of the map shows the coordinates "-13077117.73391, 4".

# Agenda



- › The ecosystem of SpatialHadoop
  - › Motivation
  - › Internal system design
  - › Applications
  - › Related work
  - › Performance Results
- › Other research projects
- › Future work

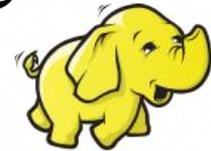
# Other Big Spatial Data Systems

Parallel

**SECONDO**

*MD*-HBase

GeoSpark



Hadoop-GIS  
*Spatial Big Data Solutions*



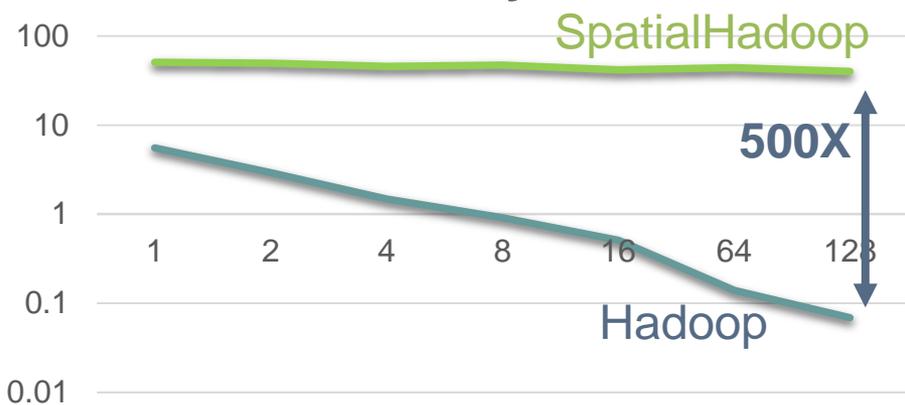
ESRI Tools  
for Hadoop

SpatialHadoop is the only extensible system that can be easily expanded by researchers and developers

# Performance Results

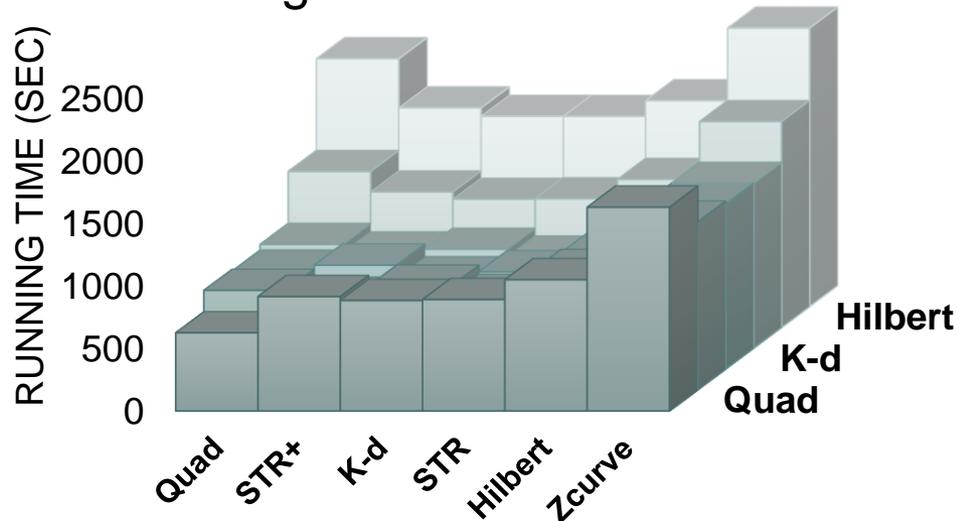


## Throughput of Range Query

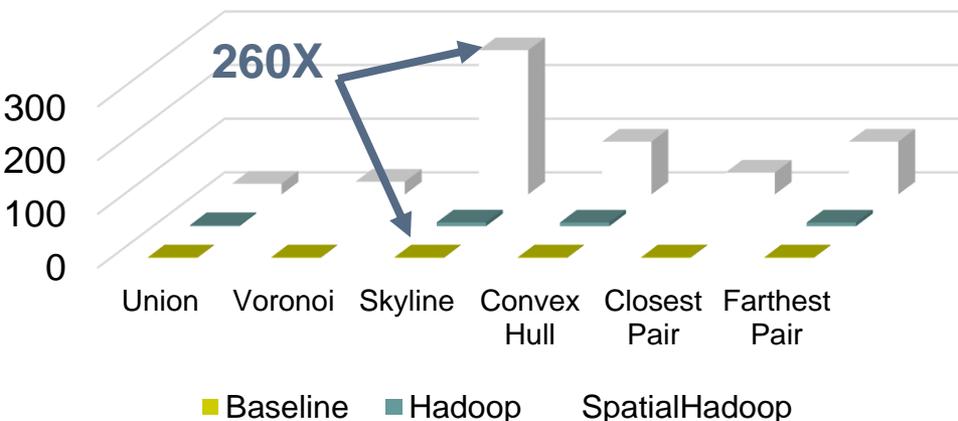


## Spatial Join

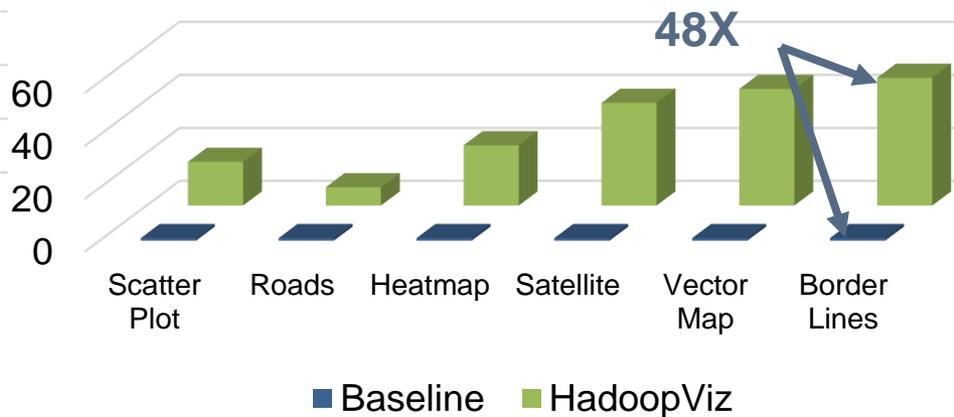
Running time with different indexes



## Speedup of CG\_Hadoop



## Visualization Speedup

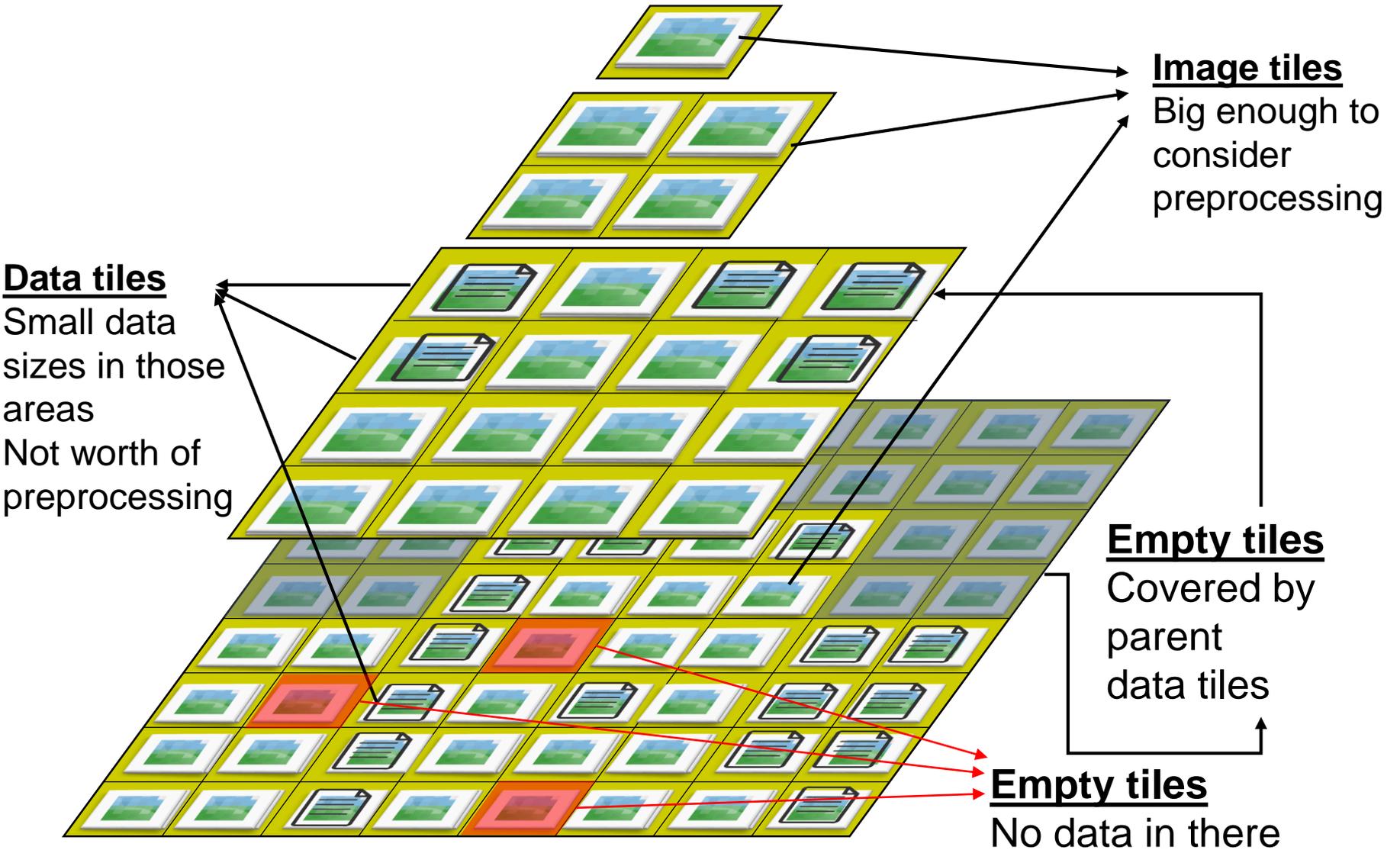


# Agenda

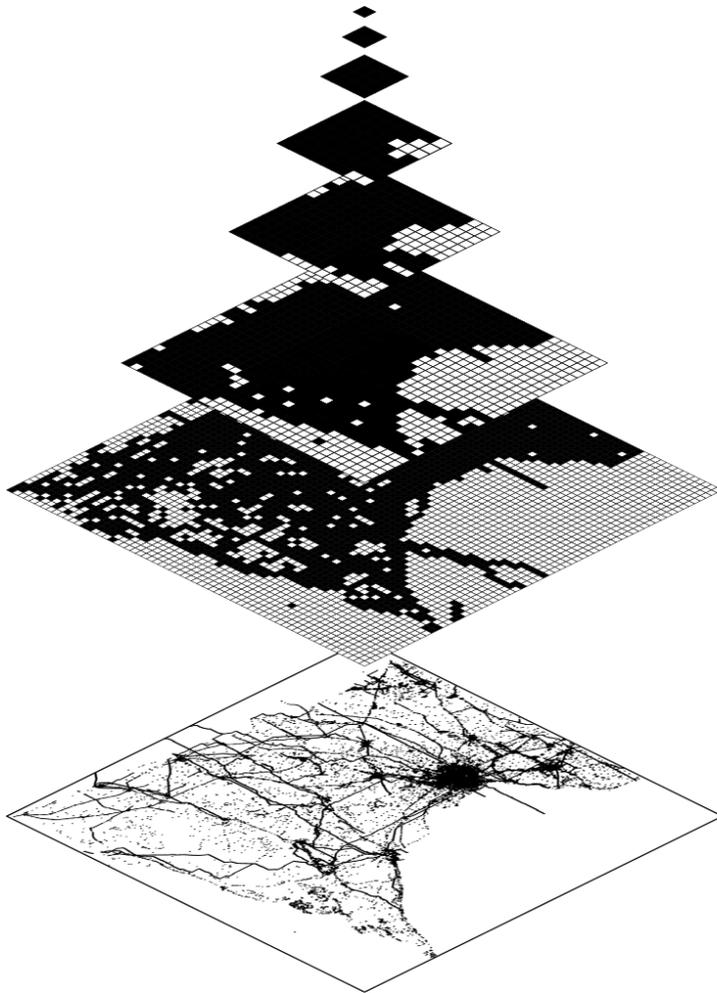


- › ~~The ecosystem of SpatialHadoop~~
  - › Motivation
  - › System design
  - › Applications
  - › Related work
  - › Performance results
  
- › Future directions

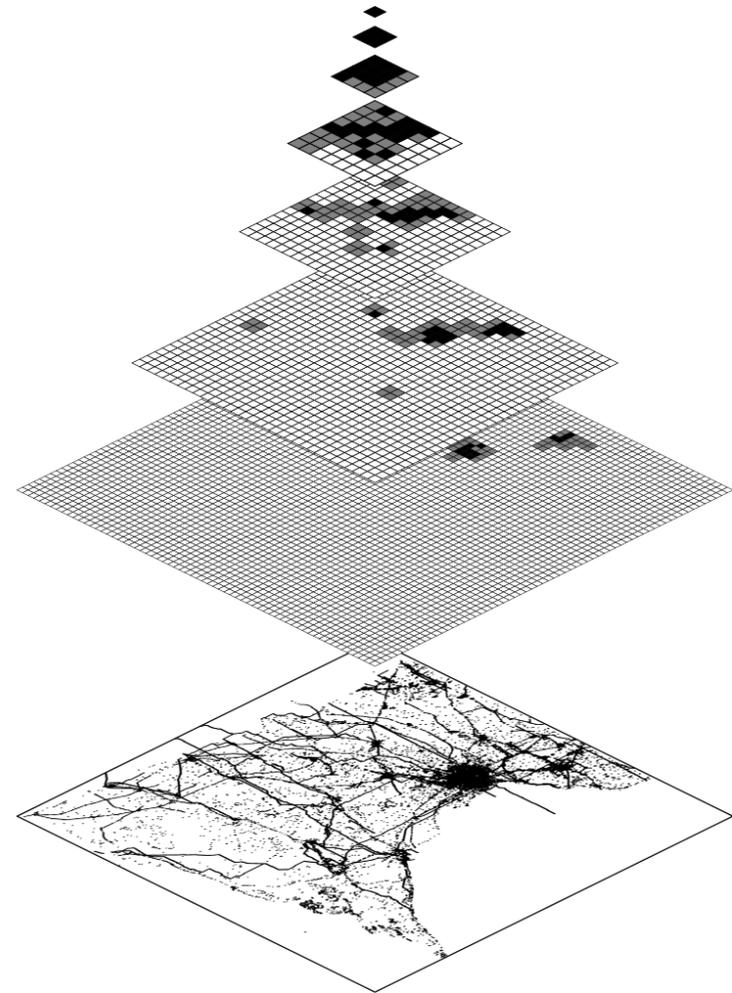
# Adaptive Multilevel Visualization



# Adaptive Multilevel Images

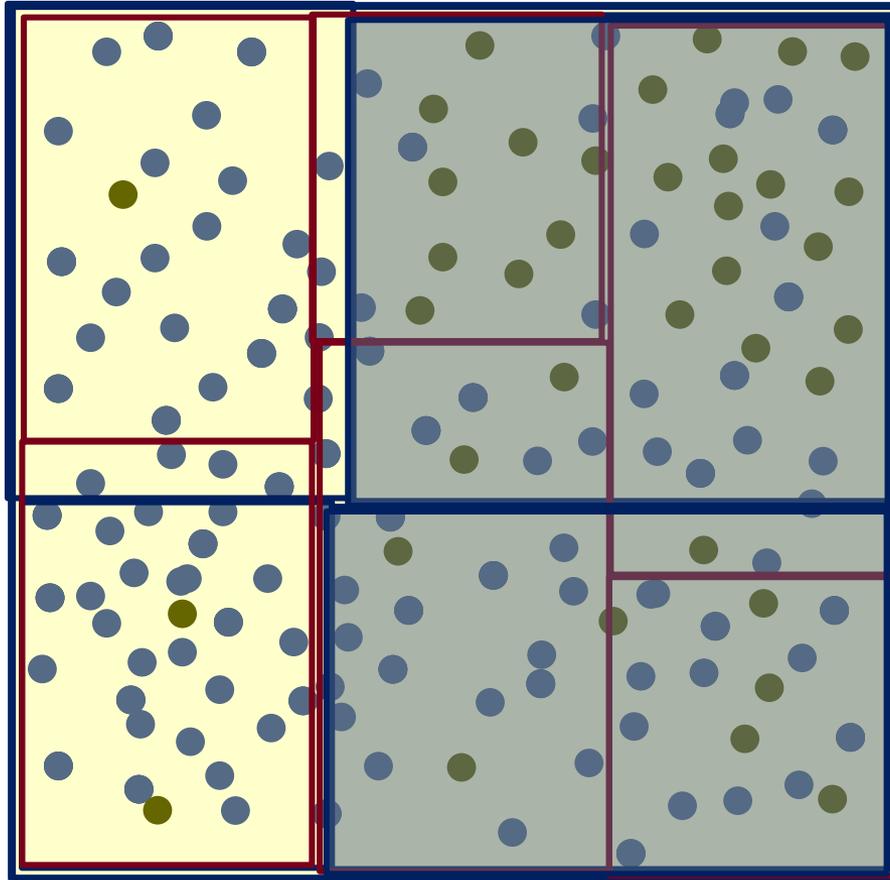


Full Image (3,160 tiles)

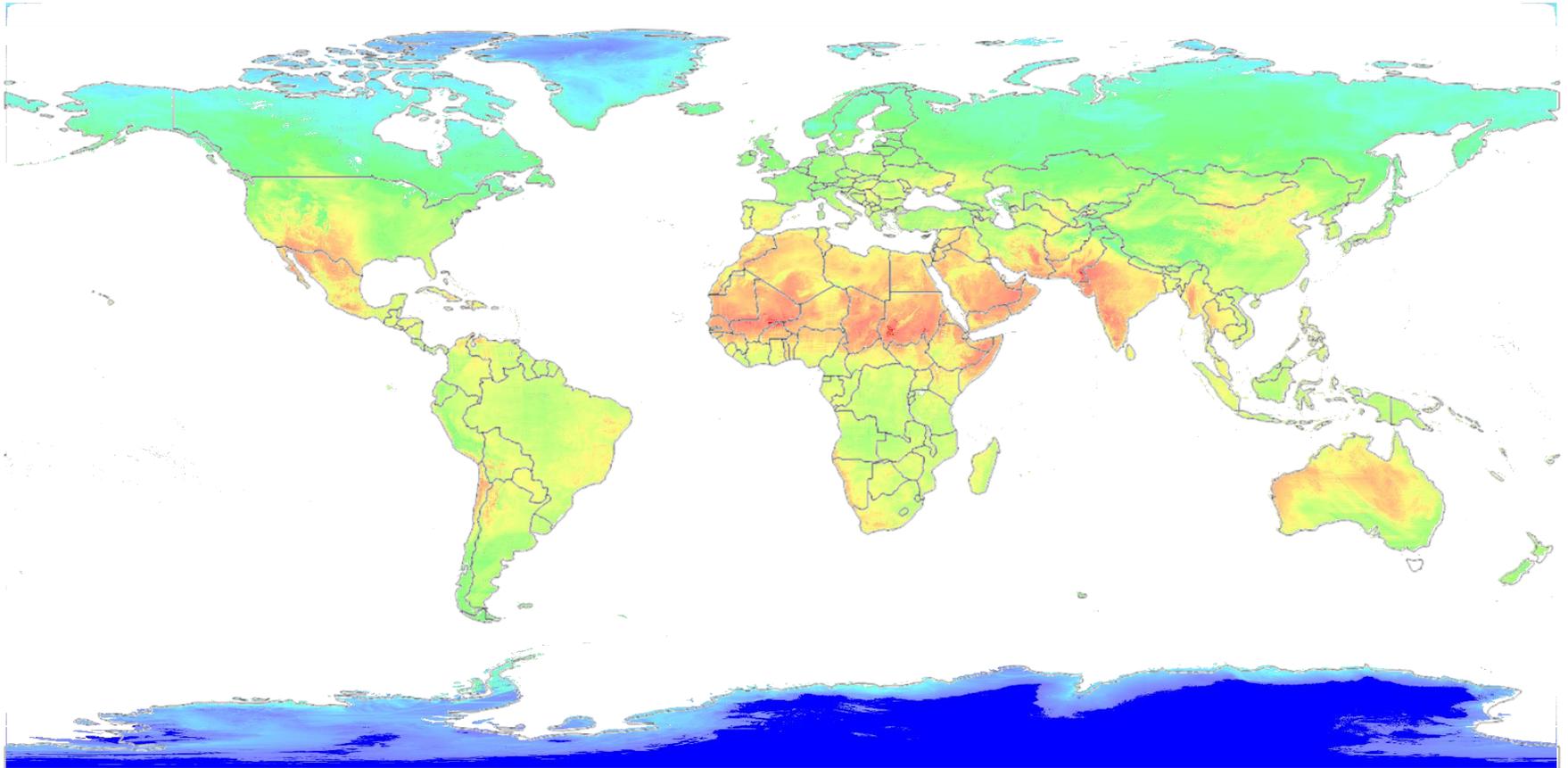


Adaptive Image (231 files)

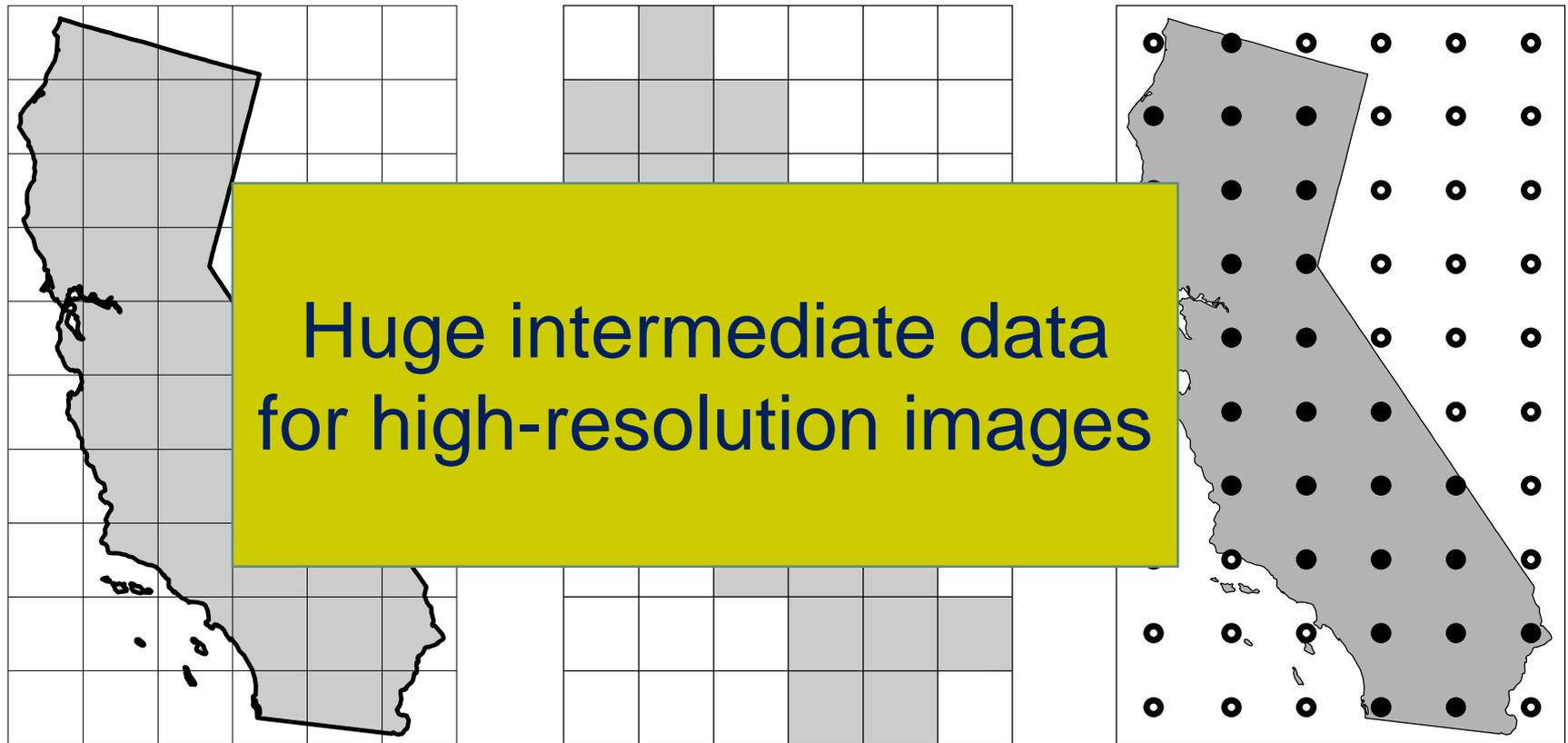
# Dynamic Indexes



# Analysis of Satellite Data



# Existing Methods



Input

Rasterize

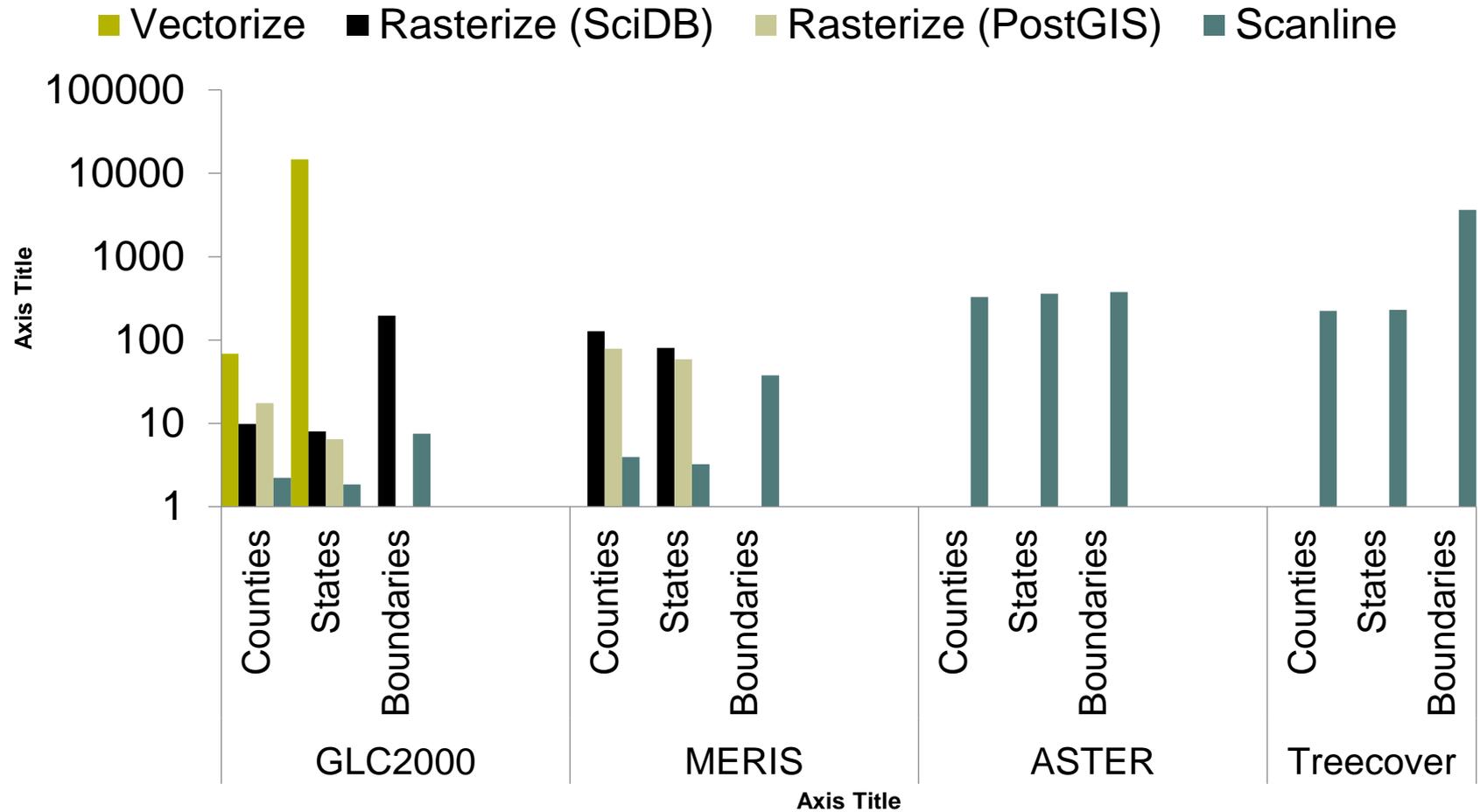
Vectorize



# Performance



Chart Title



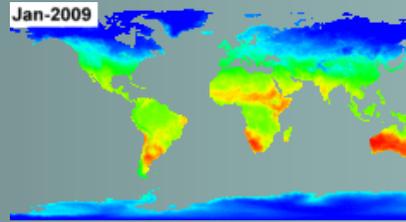
# Summary



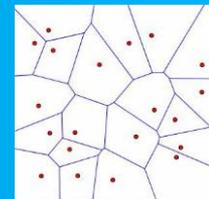
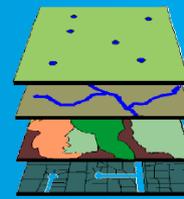
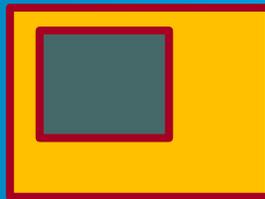
Apps



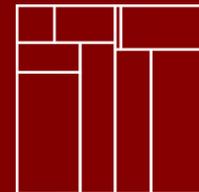
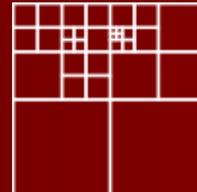
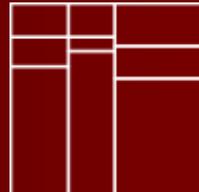
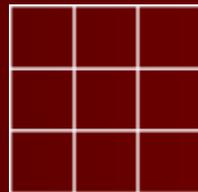
Visualization  
(HadoopViz)



Operations



Indexes



**Thank You**

Questions?