

Interactive Exploration of Big Spatial Data on SpatialHadoop and Beyond

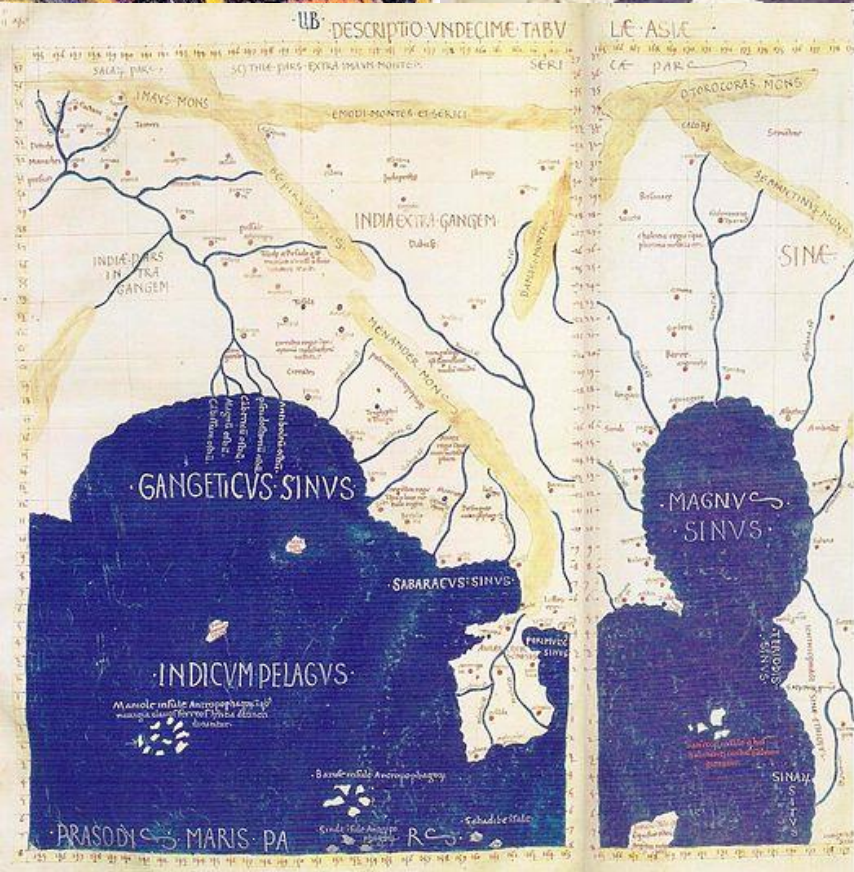
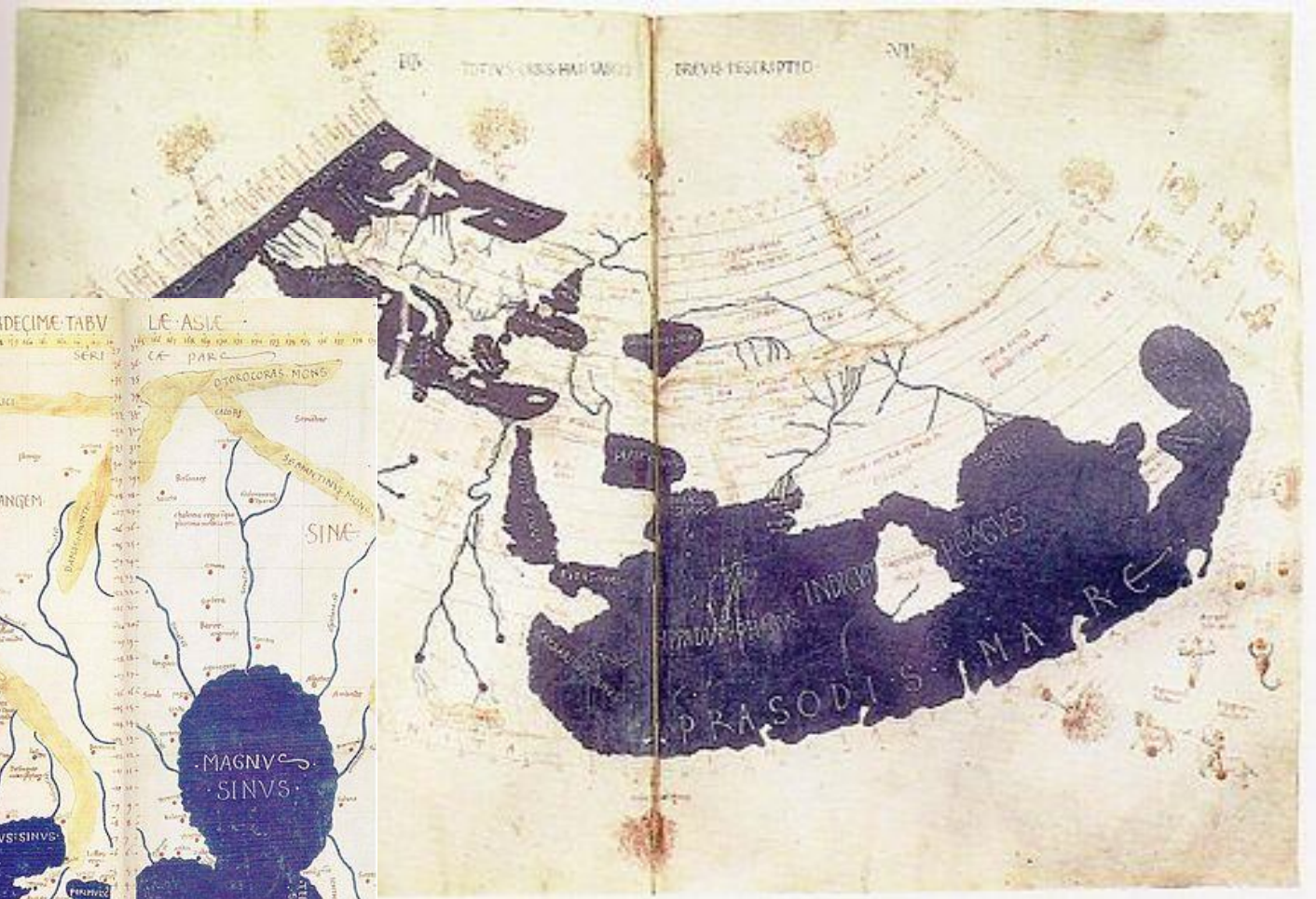
Ahmed Eldawy

Computer Science and Engineering

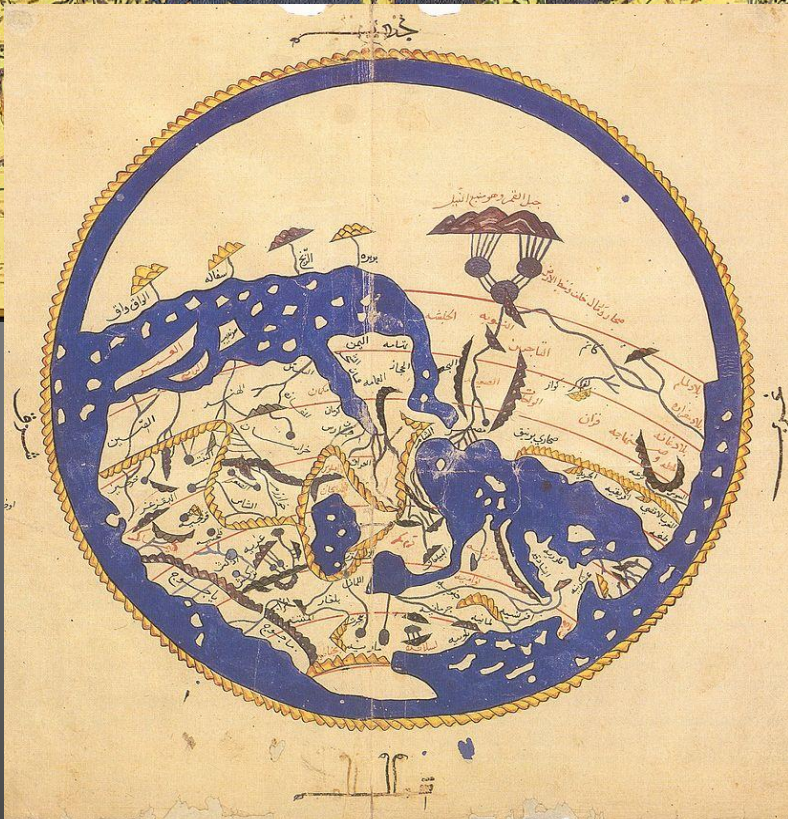
Once upon a
time...



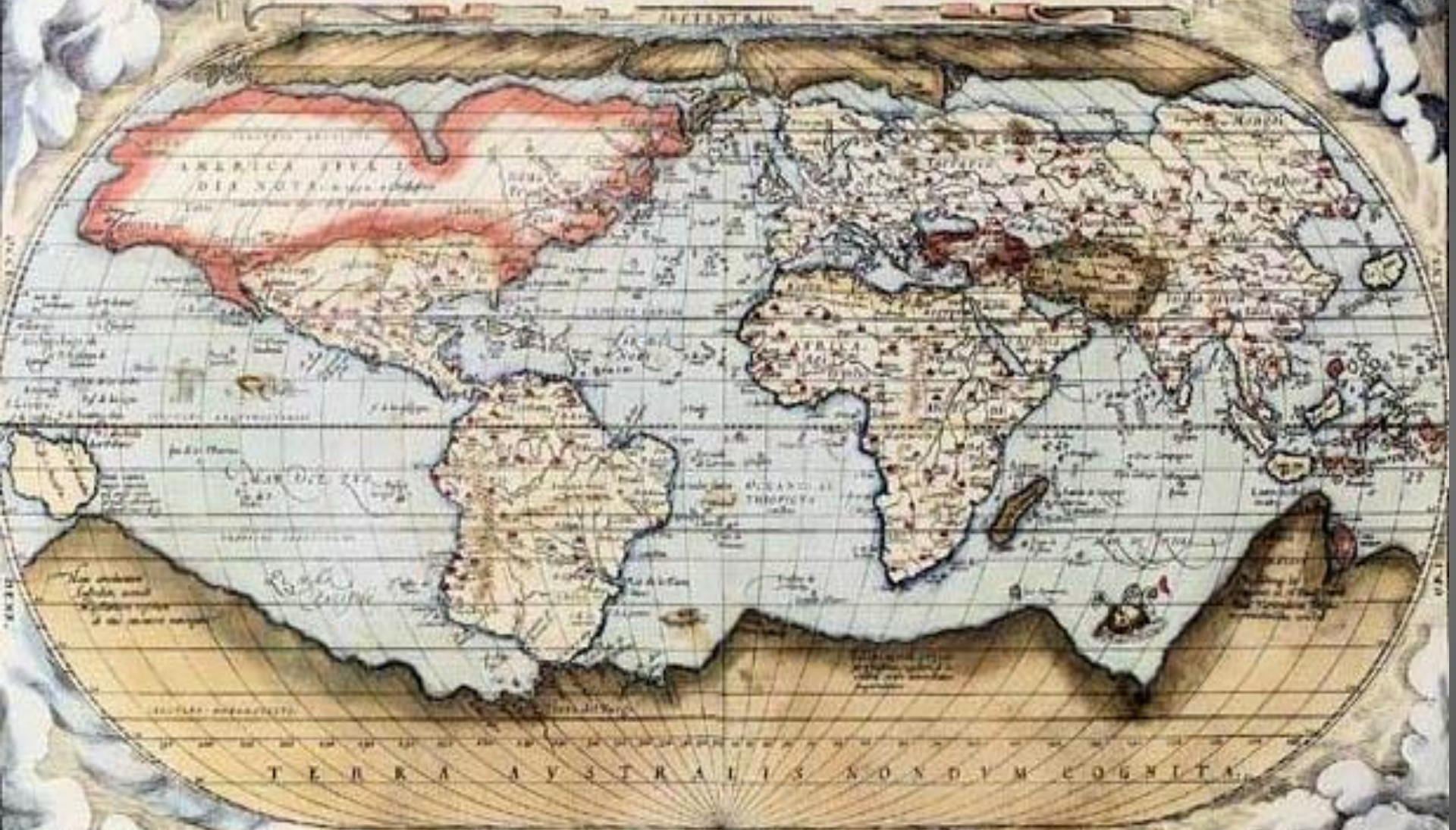
Claudius Ptolemy (AD 90 – AD 168)



Al Idrisi (1099–1165)



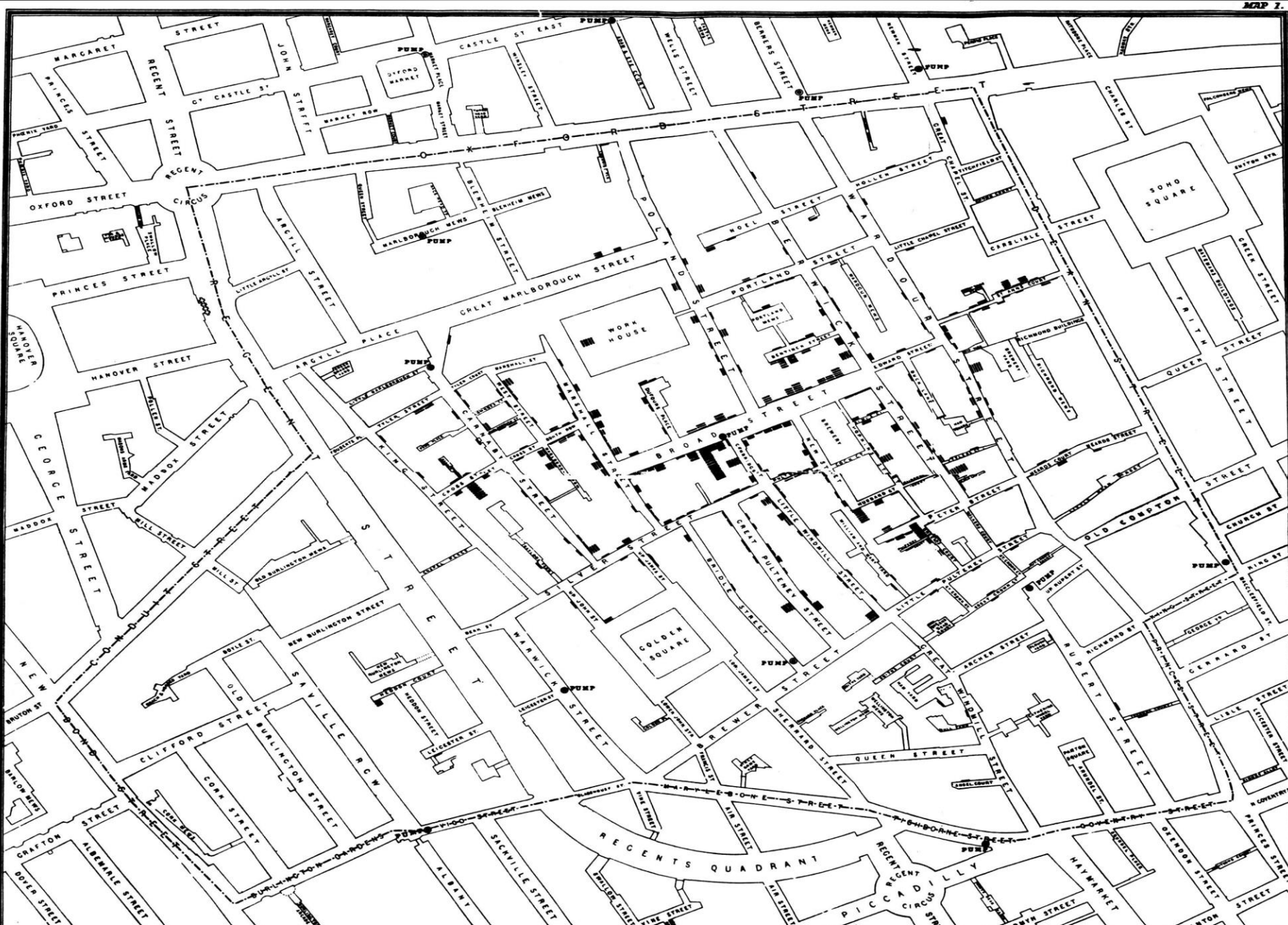
TYPVS ORBIS TERRARVM



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



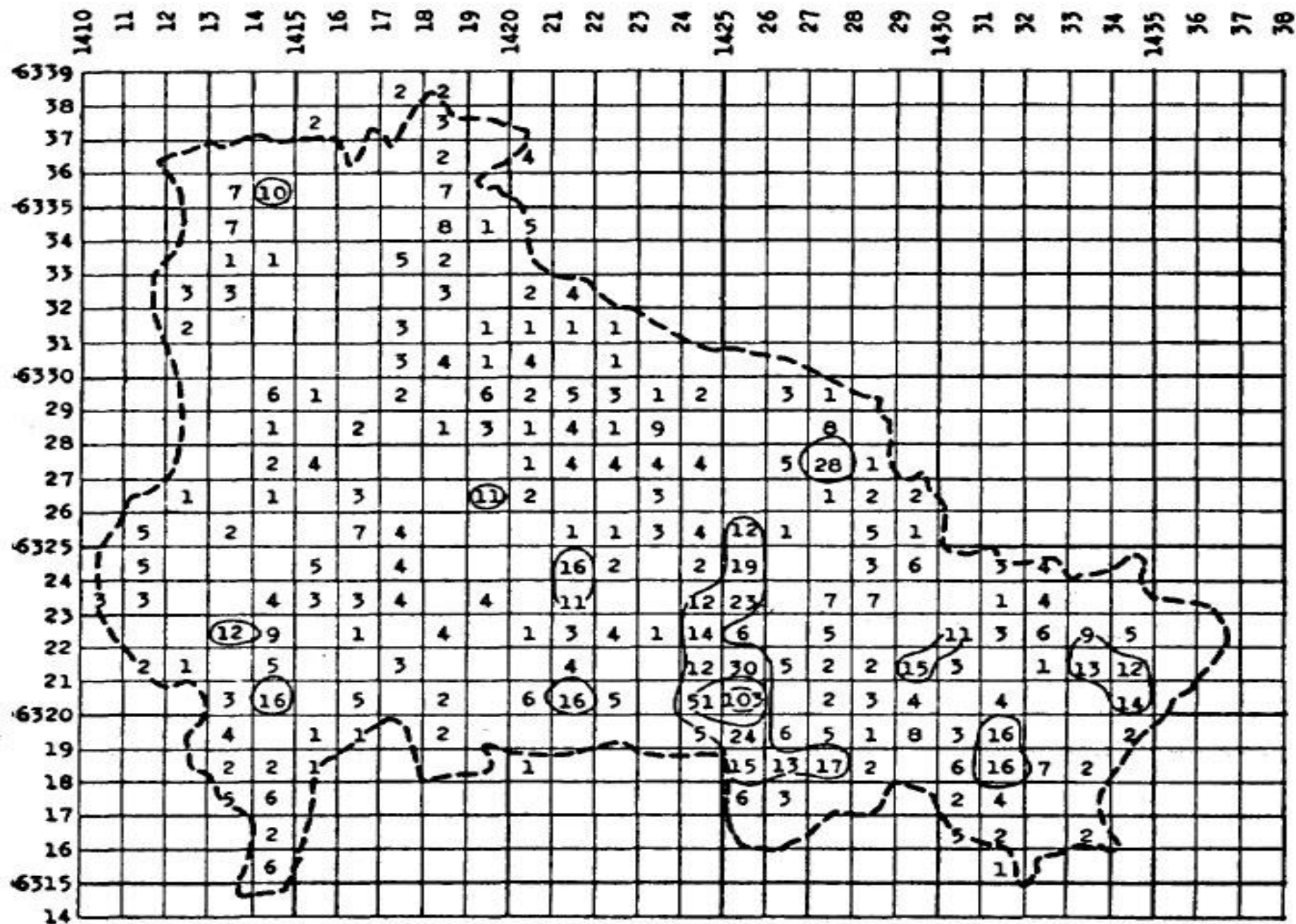
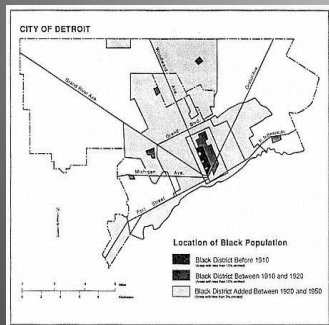
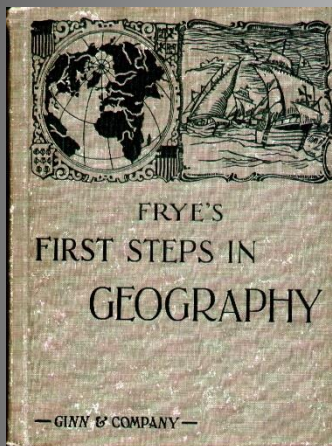
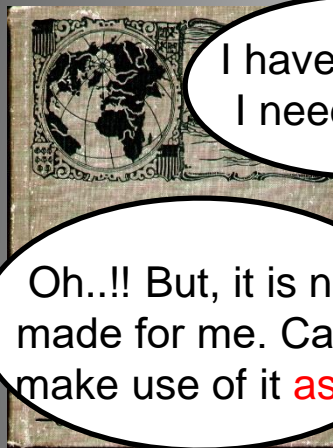


FIGURE 3—Children under 15 years of age in 1940.





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



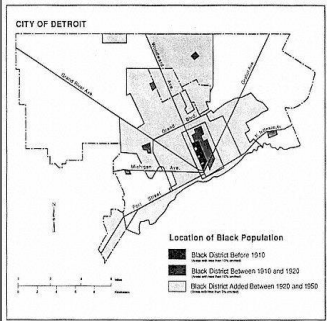
I have **BIG** data.
I need **HELP**..!!



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

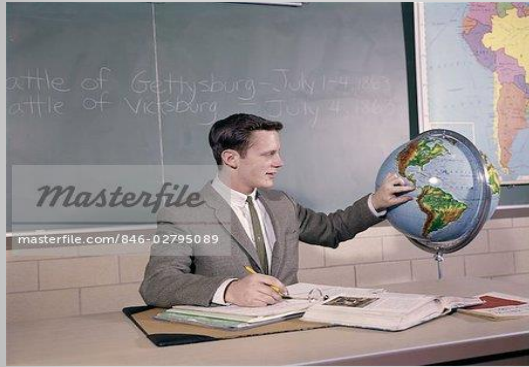


My pleasure. Here it is.



masterfile.com/846-02793618

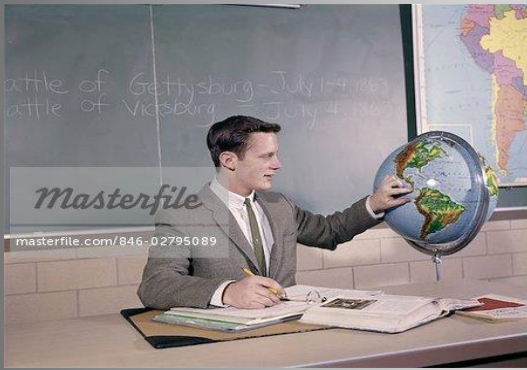




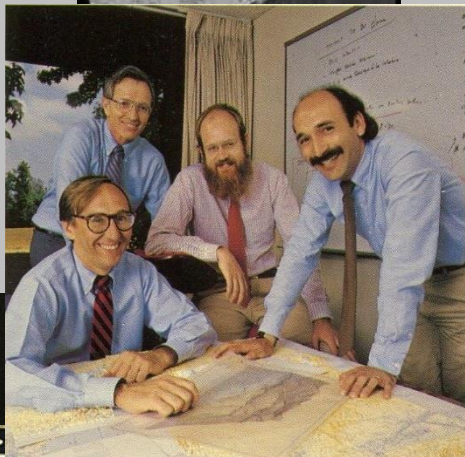
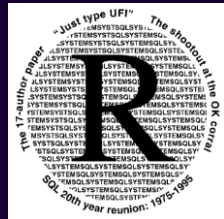
Kindly let me understand your needs

1969

Kindly let me get the technology you have



ESRI



DATABASE MANAGEMENT SYSTEMS



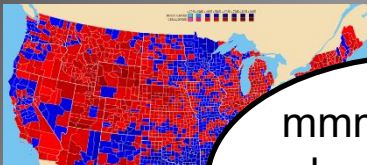
Informix

SQL

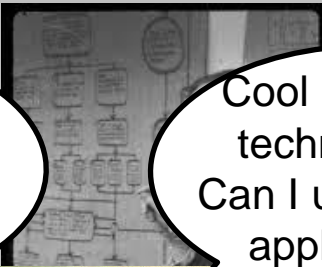


ESRI

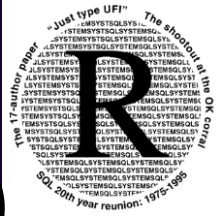




mmm...Let me check with my good friends there.



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



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**

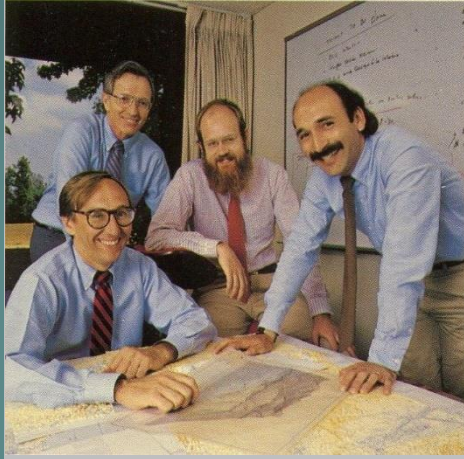
My pleasure. Here it is.



Informix

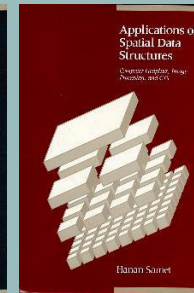
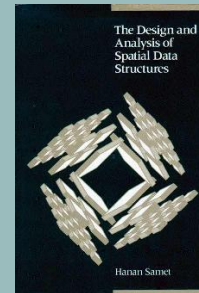
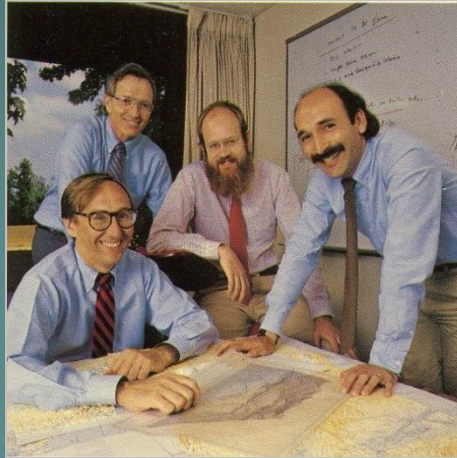
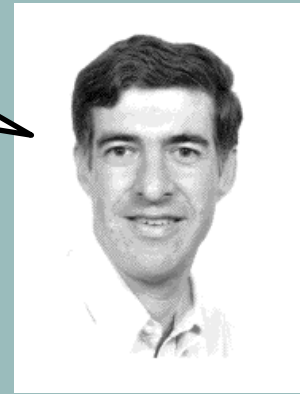
SQL

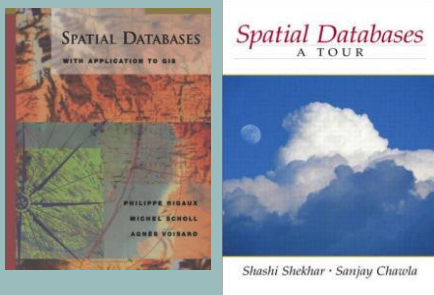
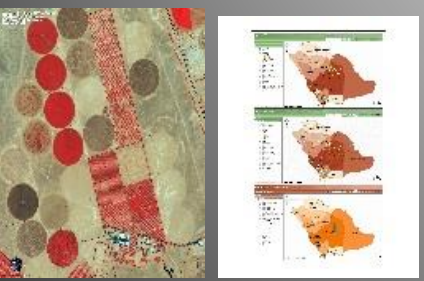
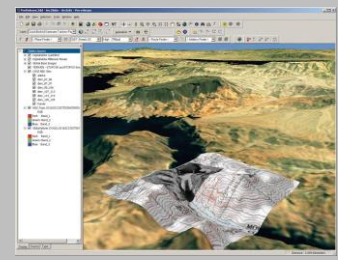
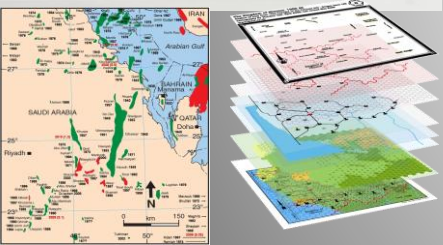
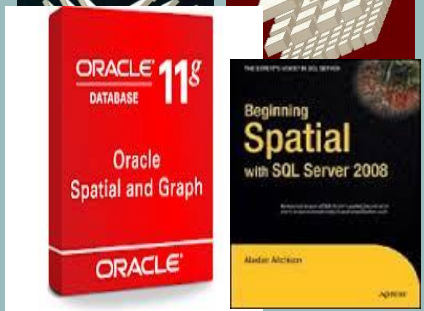
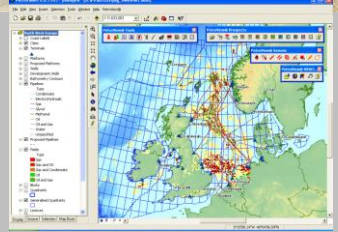
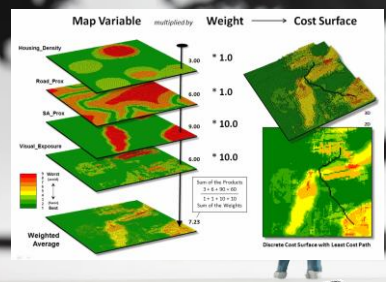


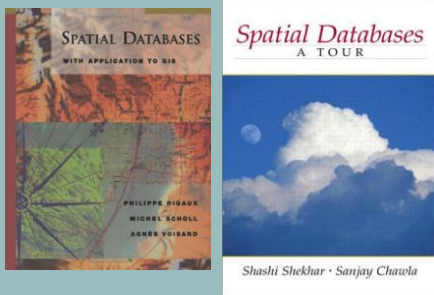
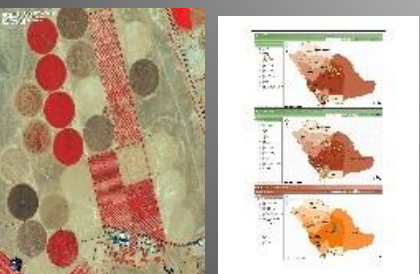
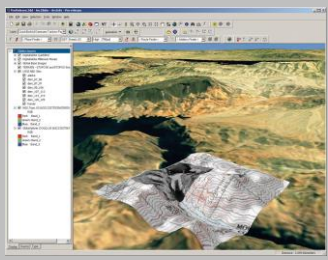
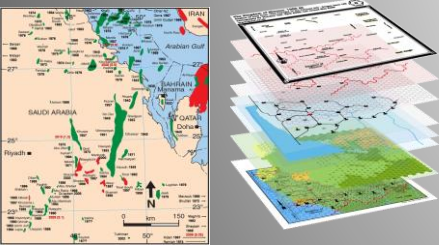
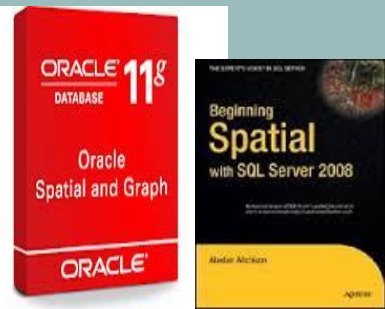
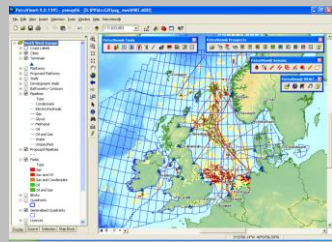
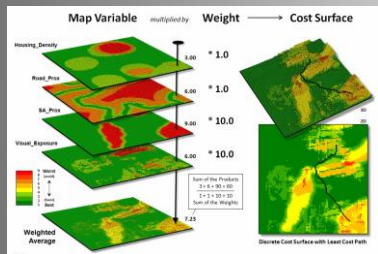


Kindly let me understand your needs

Kindly let me get the technology you have



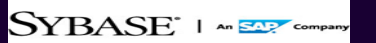
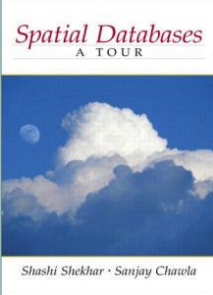
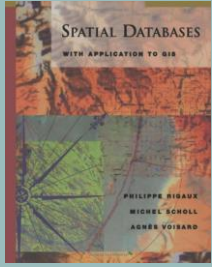
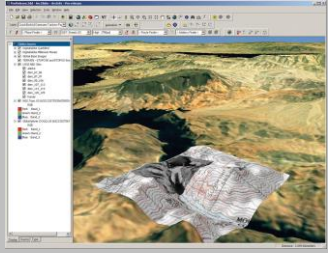
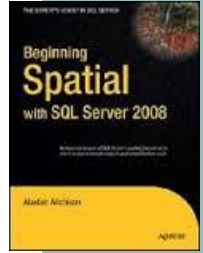
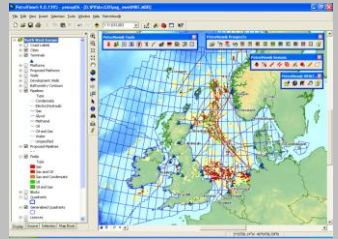






facebook

MapReduce



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

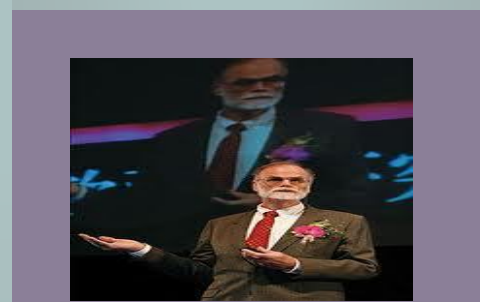
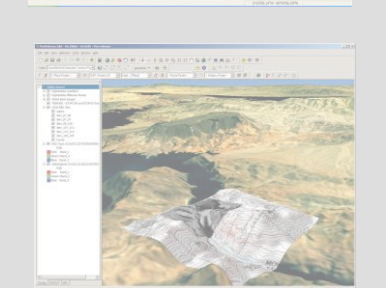
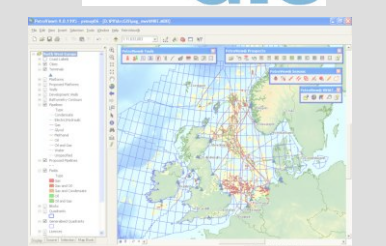
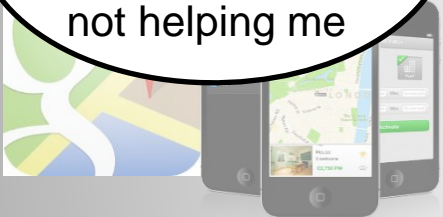
Cool **Big Data** technology..!!
Can I use it in my application?

My pleasure.
Here it is.

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

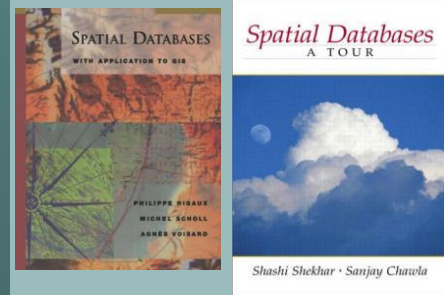
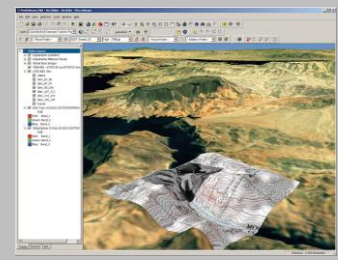
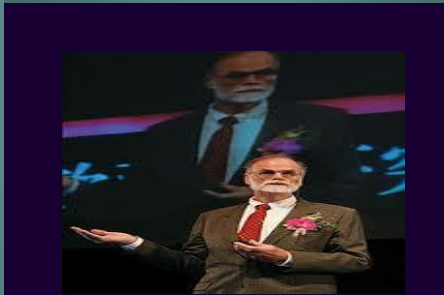
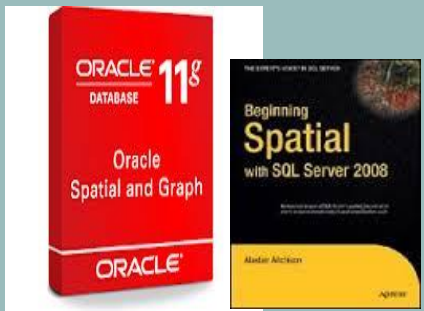
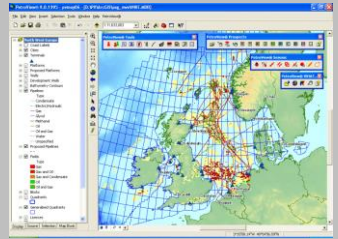
Sorry, seems like
the DBMS
technology cannot
scale more

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





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



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

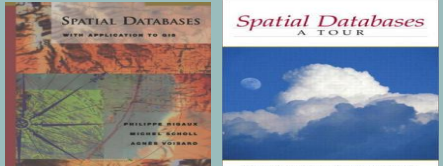
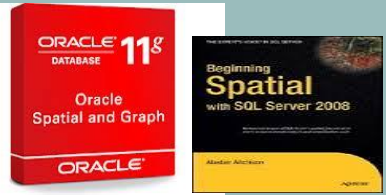
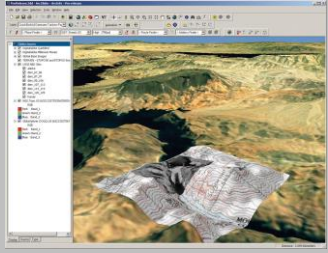
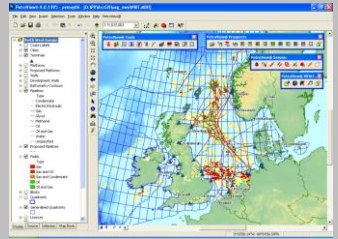




Kindly let me understand your needs



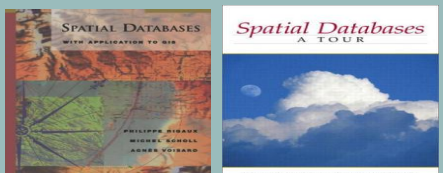
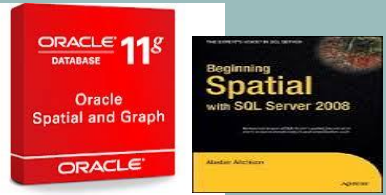
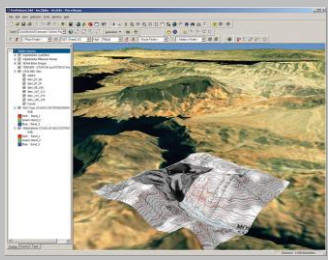
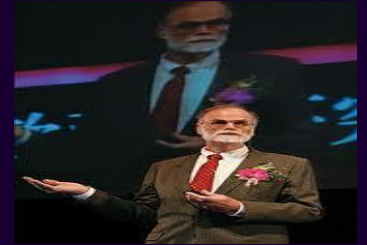
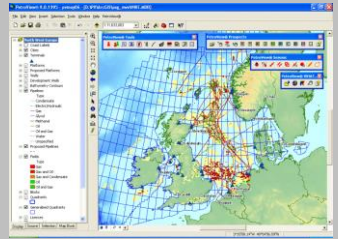
Kindly let me get the technology you have





Big Spatial Data

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



Tons of Spatial data out there...

twitter



Geotagged Microblogs



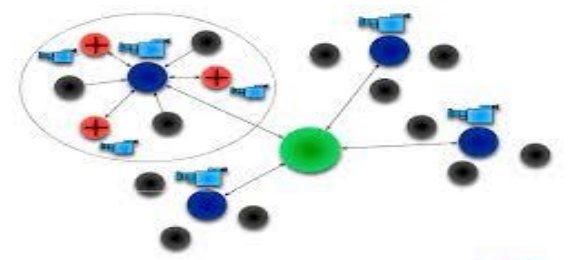
Geotagged Pictures



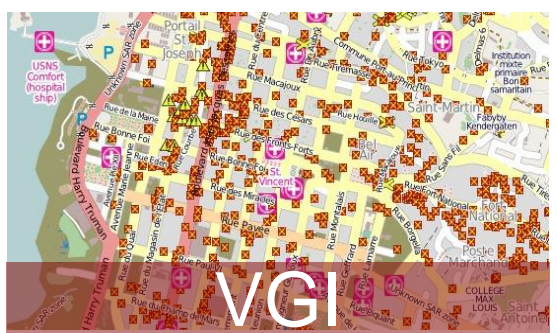
Medical Data



Smart Phones



Sensor Networks



VGI

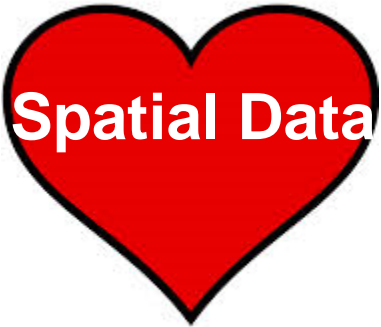


Satellite Images



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;
```

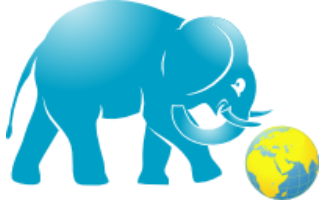
```
points = LOAD 'points' AS
  (id:int, location:point);
result = FILTER points BY
  Overlap(location, rectangle
  (xmin, ymin, xmax, ymax));
```



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



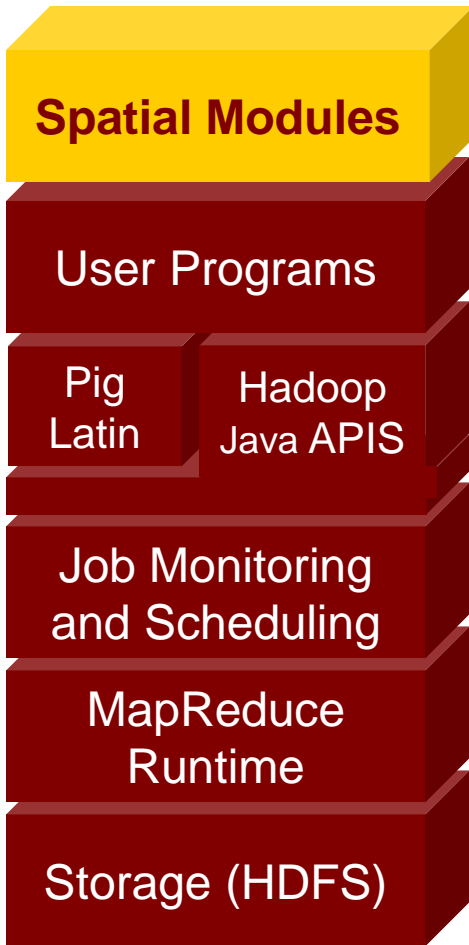
Collaboration



>500GB public datasets for
benchmarking and testing

The Built-in Approach of SpatialHadoop

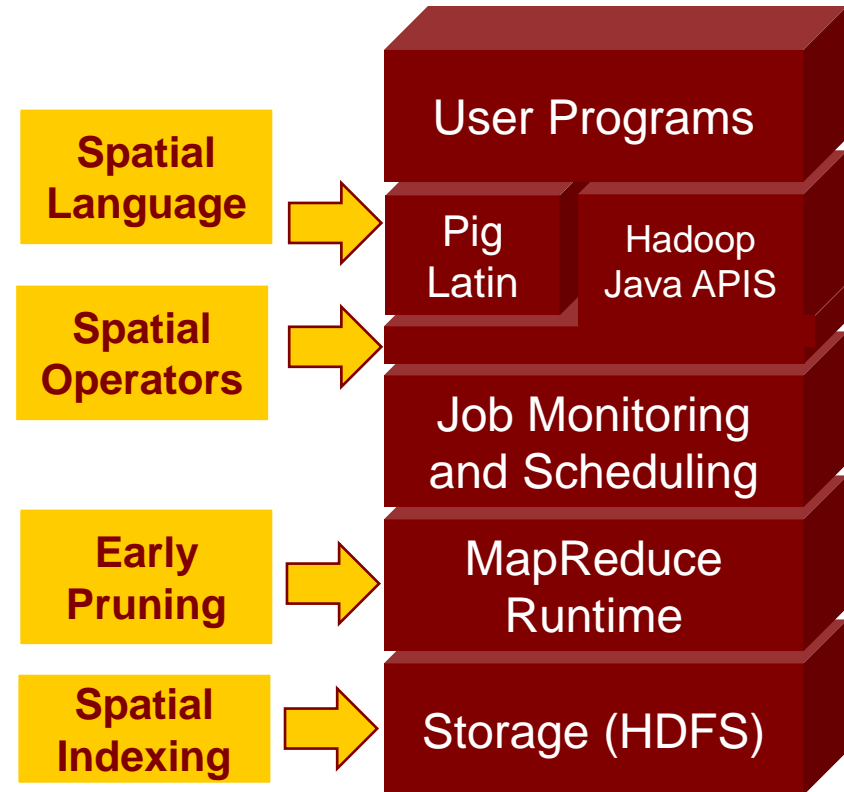
The On-top Approach



From Scratch Approach



The Built-in Approach (SpatialHadoop)



Agenda

- The ecosystem of SpatialHadoop
 - Motivation
 - Internal system design
 - Applications
 - Related work
 - Performance results
- Interactive data exploration

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

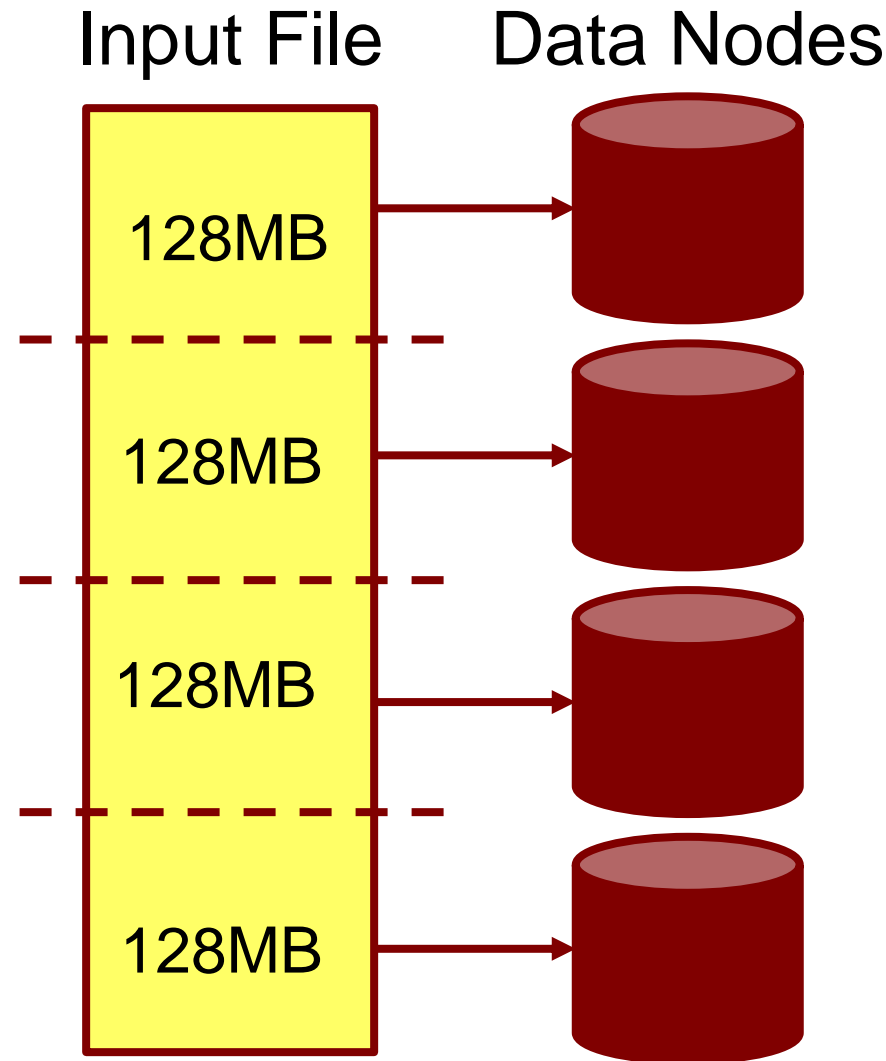
Indexing

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

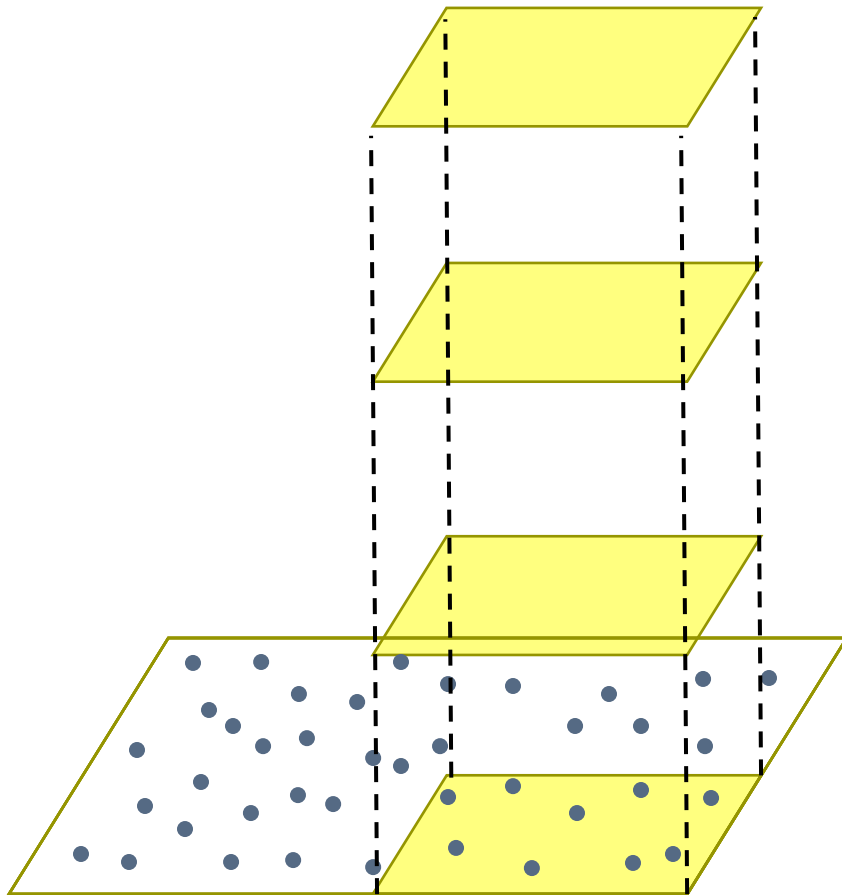
ST-Hadoop

Data Loading in Hadoop

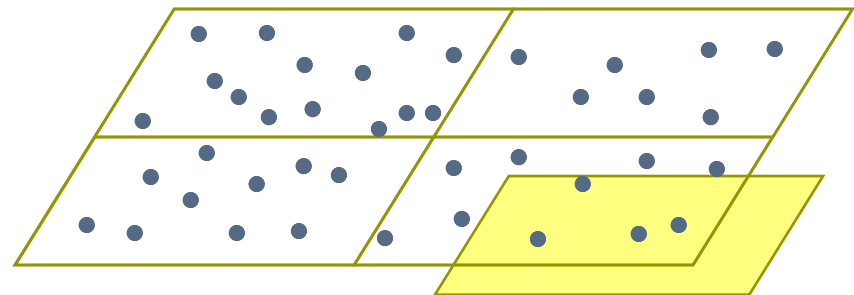
- ▶ Blindly chops down a big file into 128MB chunks
- ▶ Values of records are not considered
- ▶ Relevant records are typically assigned to two different blocks
- ▶ HDFS is too restrictive where files cannot be modified



Spatial Distributed File System

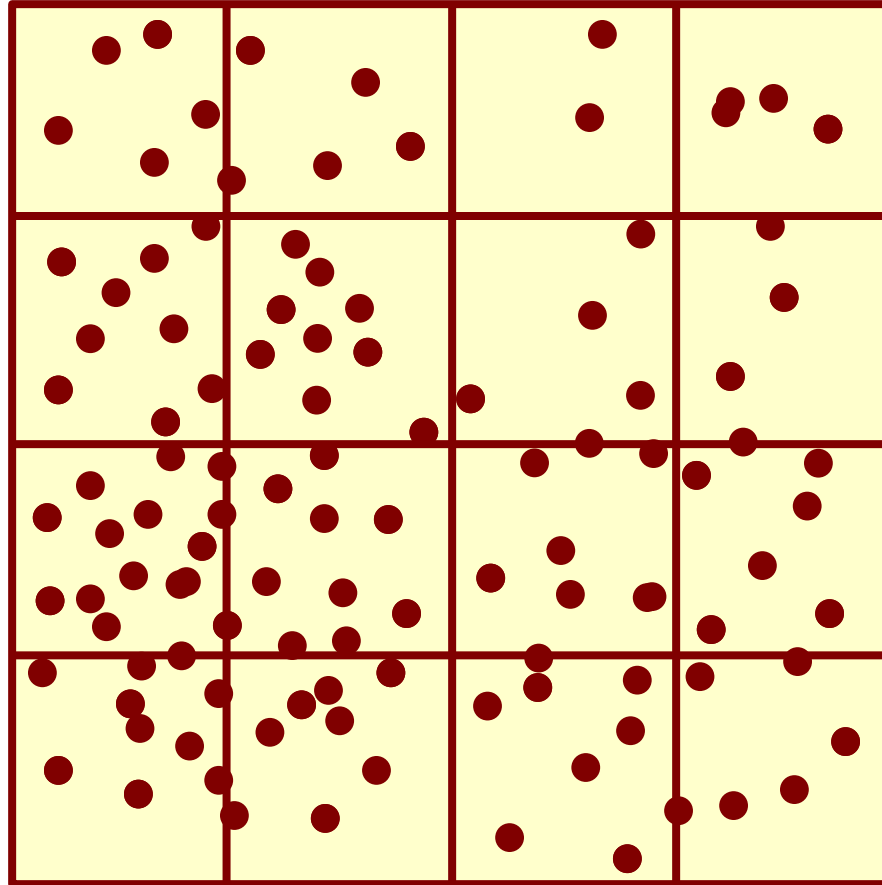


4/8/2021 **Default Partitioning**



Spatial Partitioning

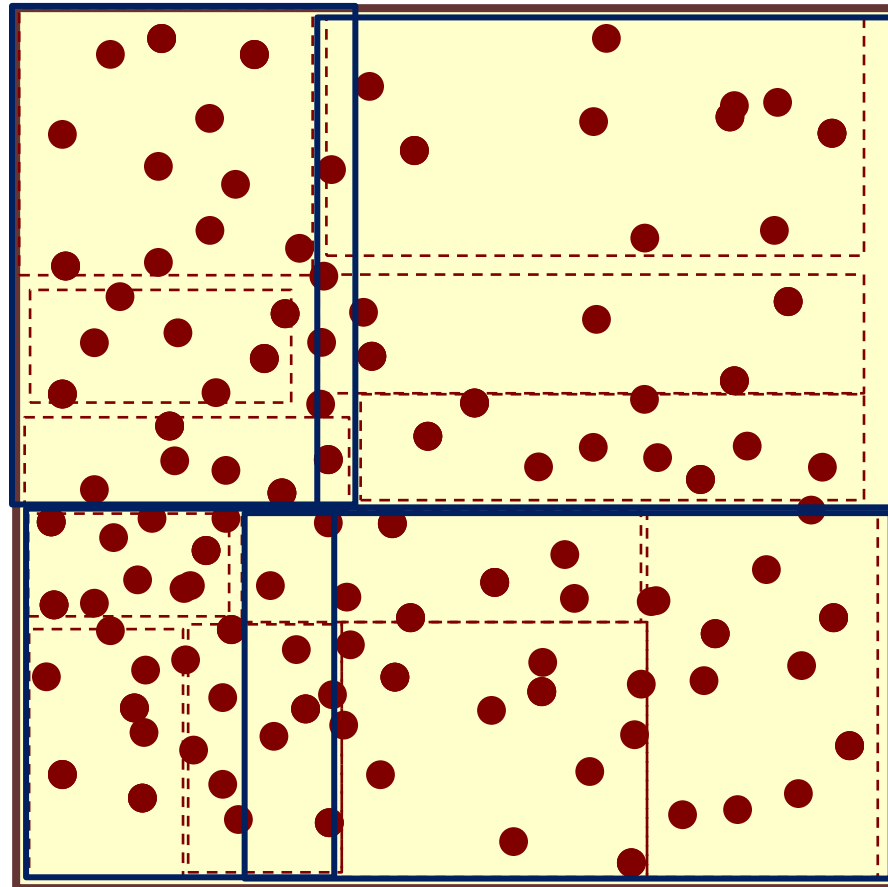
Uniform Grid



Works only for uniformly distributed data

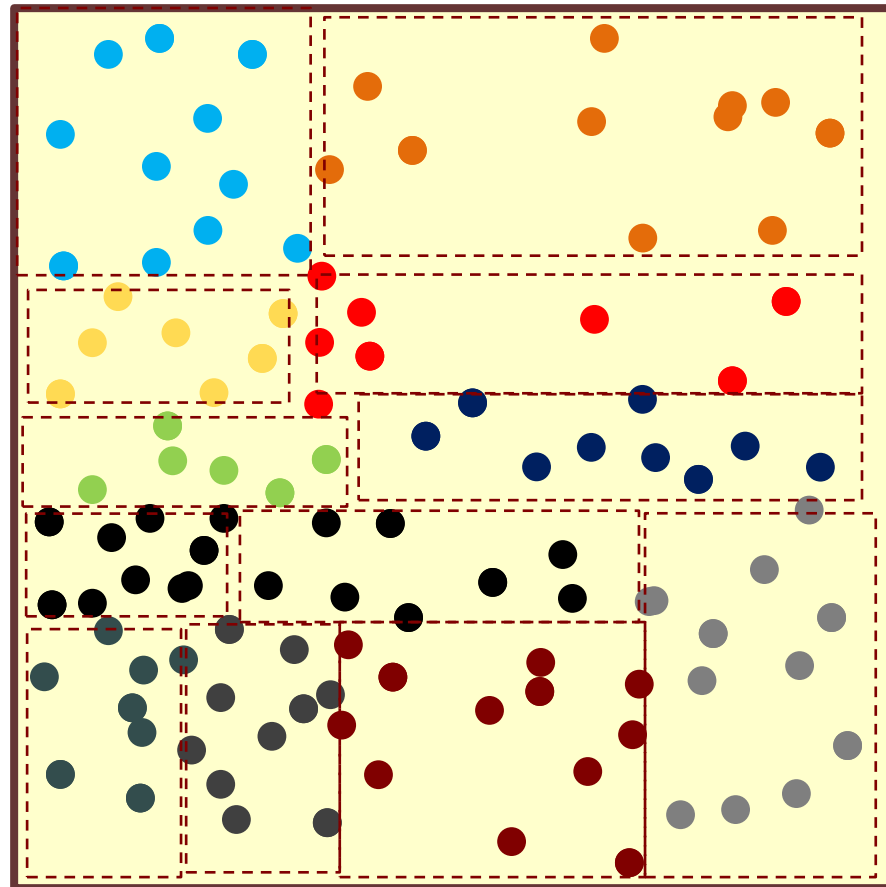
R-tree

- › 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
 - › α : 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
 - › α : Index overhead
- › Use MBR of leaf nodes as partition boundaries
- › Partition the data



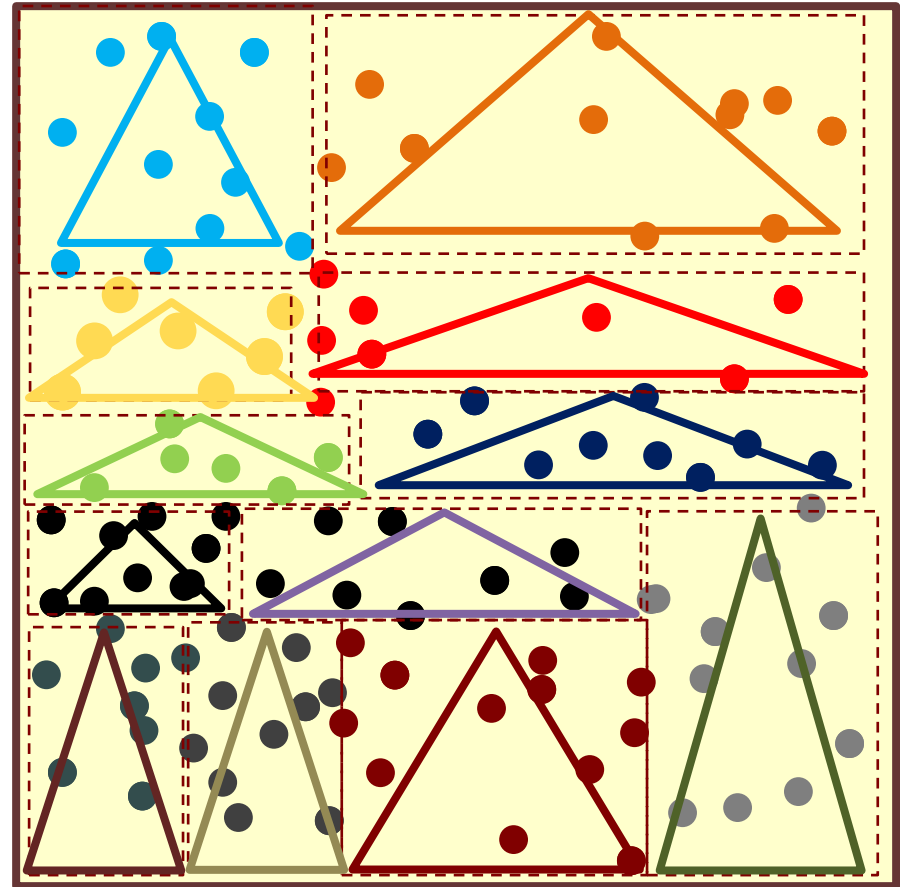
R-tree

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

- › Leaf node capacity C

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

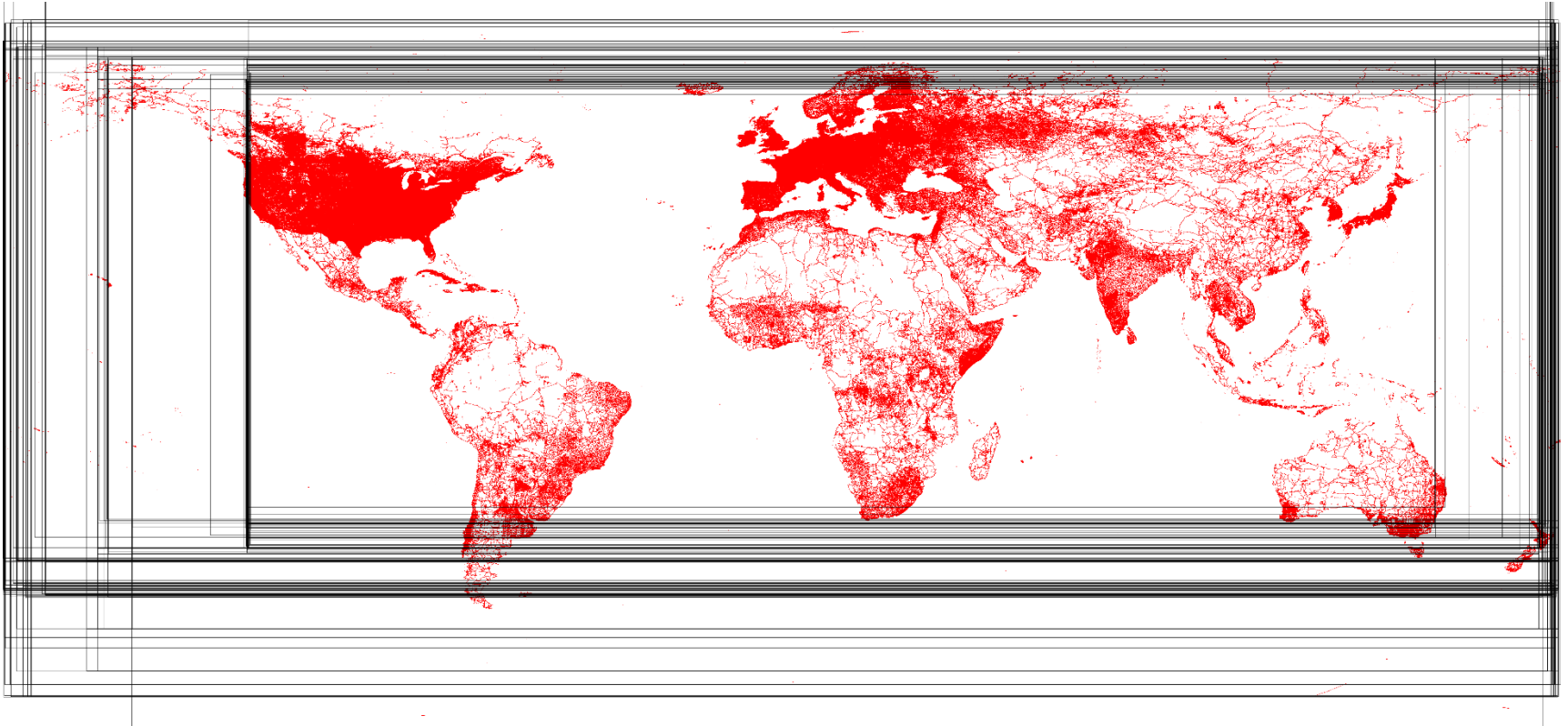
- › k : Sample size
- › B : HDFS Block capacity
- › $|R|$: Input size
- › α : Index overhead
- › Use MBR of leaf nodes as partition boundaries
- › Partition the data
- › Optional: Build R-tree Local indexes



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]

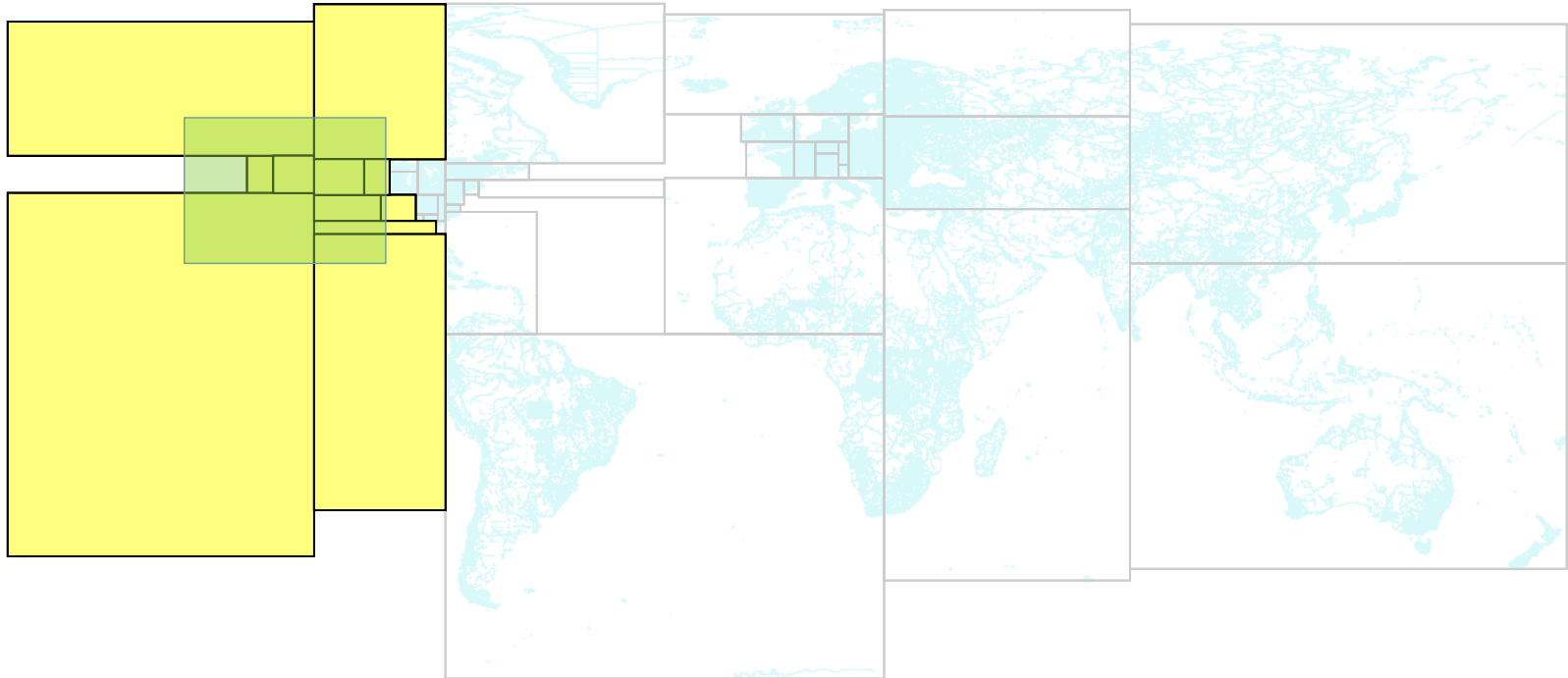
ST-Hadoop

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

4/8/2021

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

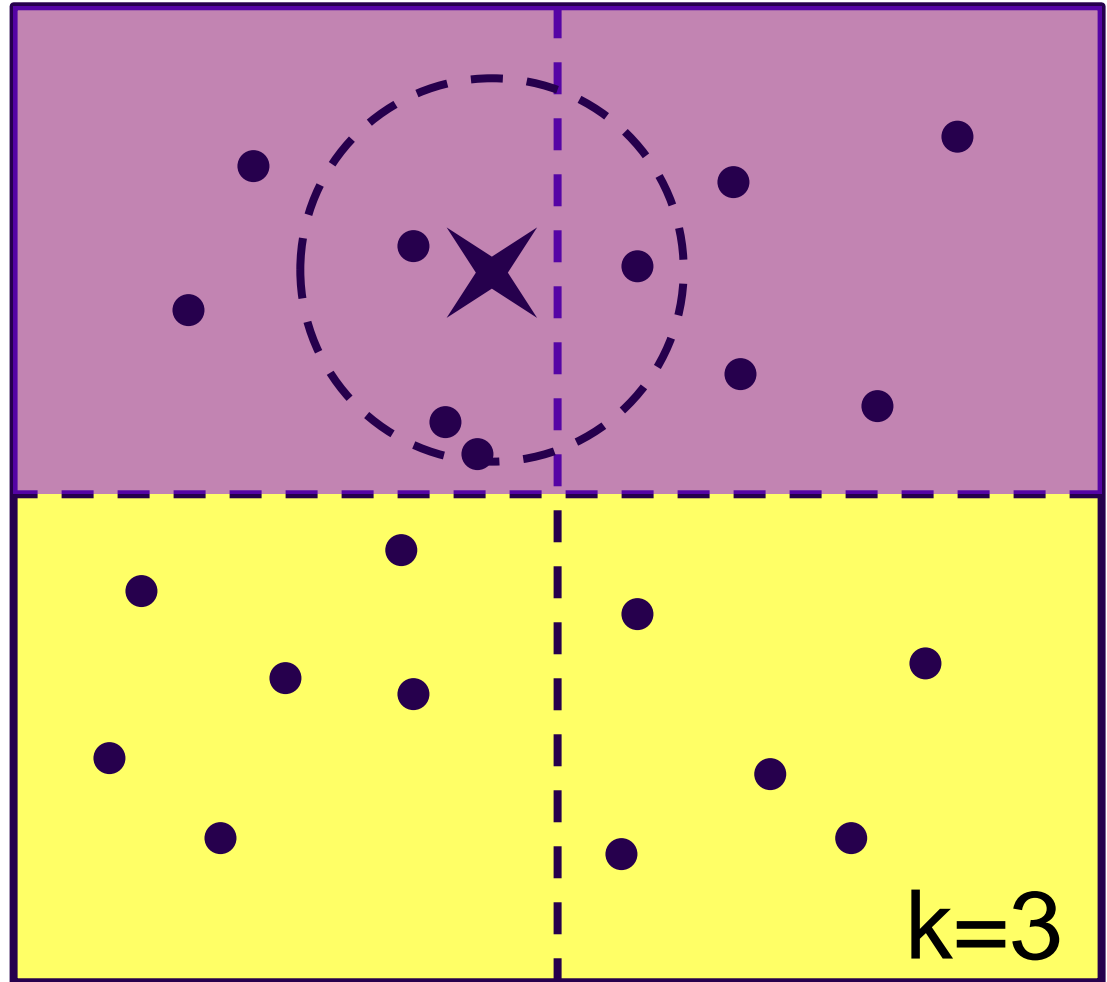
KNN over Indexed Data

First iteration runs as before and result is tested for correctness

✗ Answer is incorrect

Second iteration processes other blocks that might contain an answer

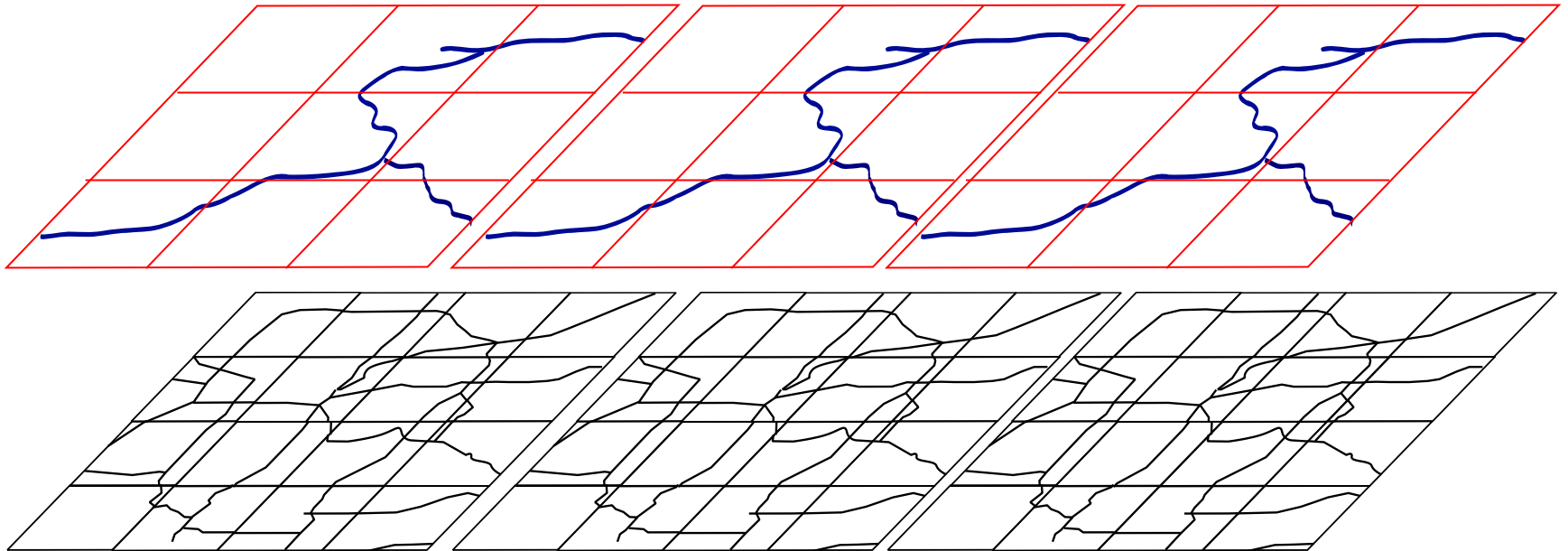
✓ Answer is correct



Spatial Join

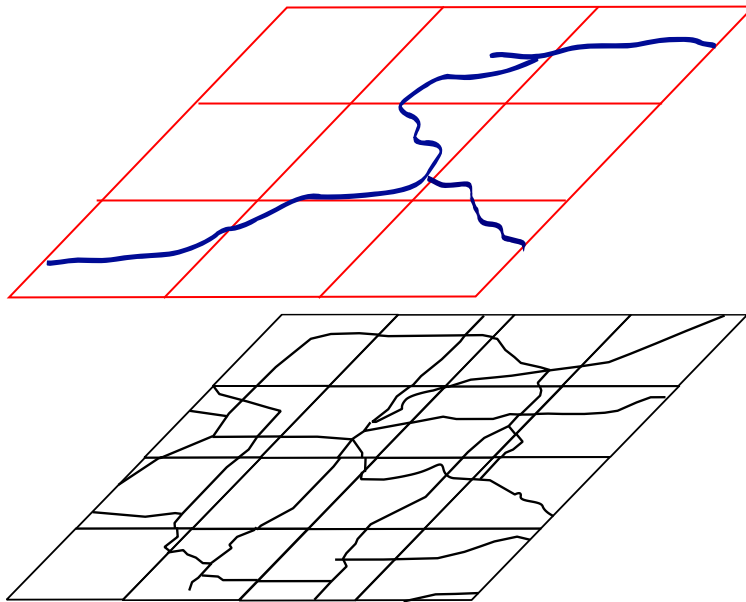
Join Directly

Partition – Join



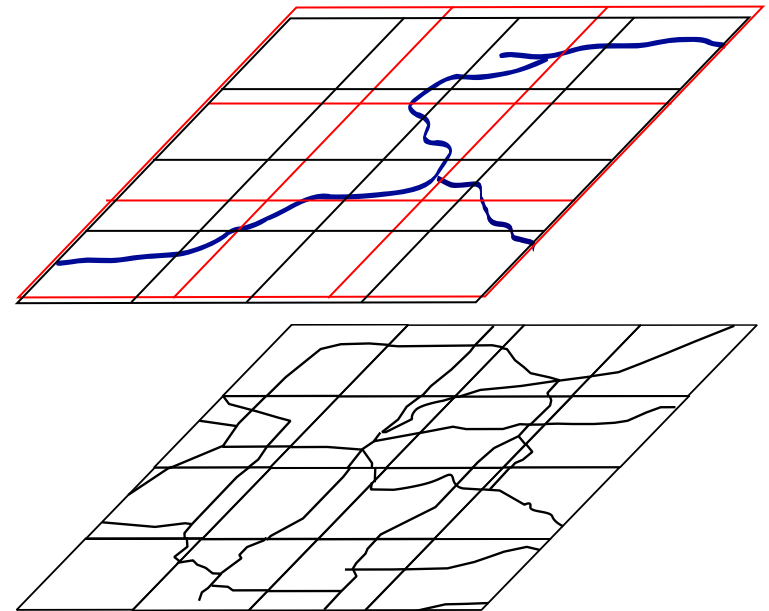
Spatial Join

Join Directly



Total of 36 overlapping pairs

Partition – Join



Only 16 overlapping pairs

CG_Hadoop

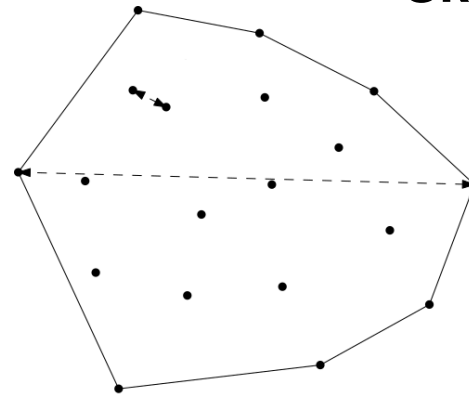
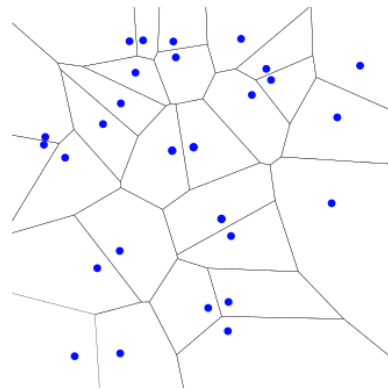
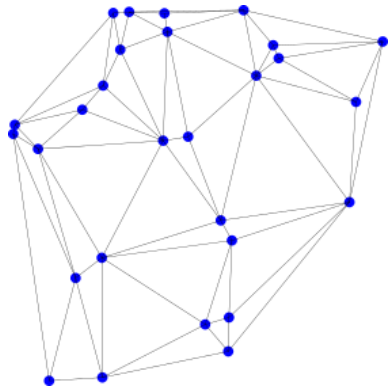


260x



Polygon Union

Skyline



**Delaunay
Triangulation**

**Voronoi
Diagram**

**Convex Hull
Farthest/closest pair**

1x

Single
Machine

29x



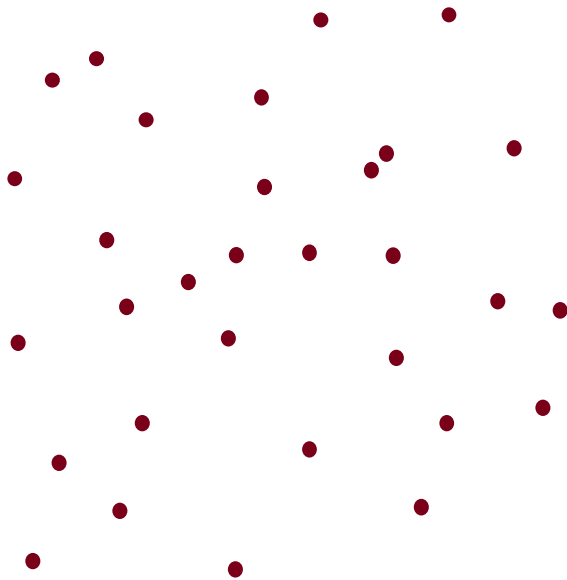
Hadoop

Spatial
Hadoop

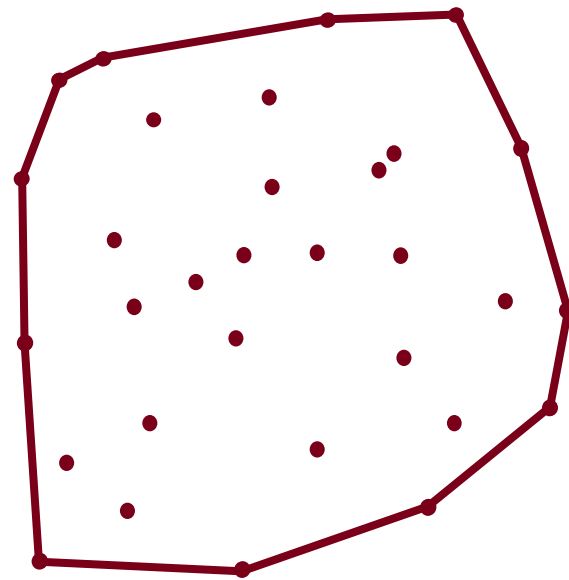
Convex Hull

Find the minimal convex polygon that contains all points

Input



Output

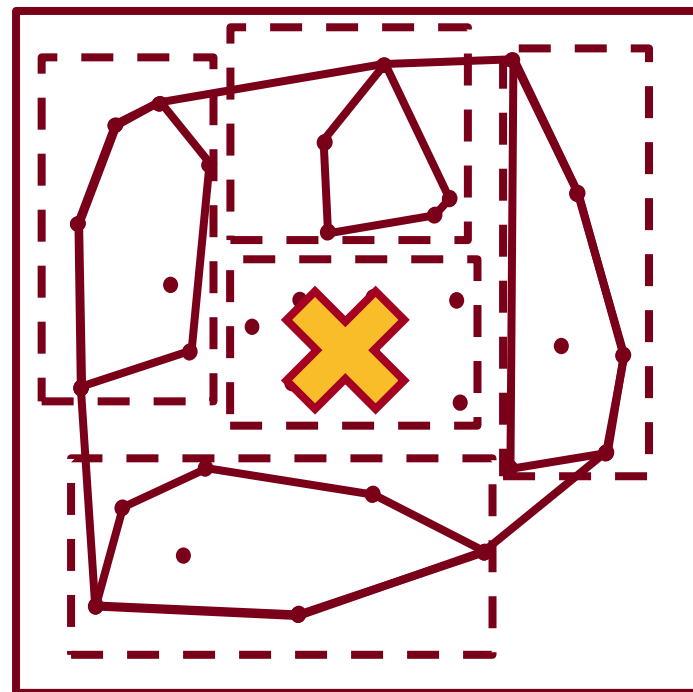
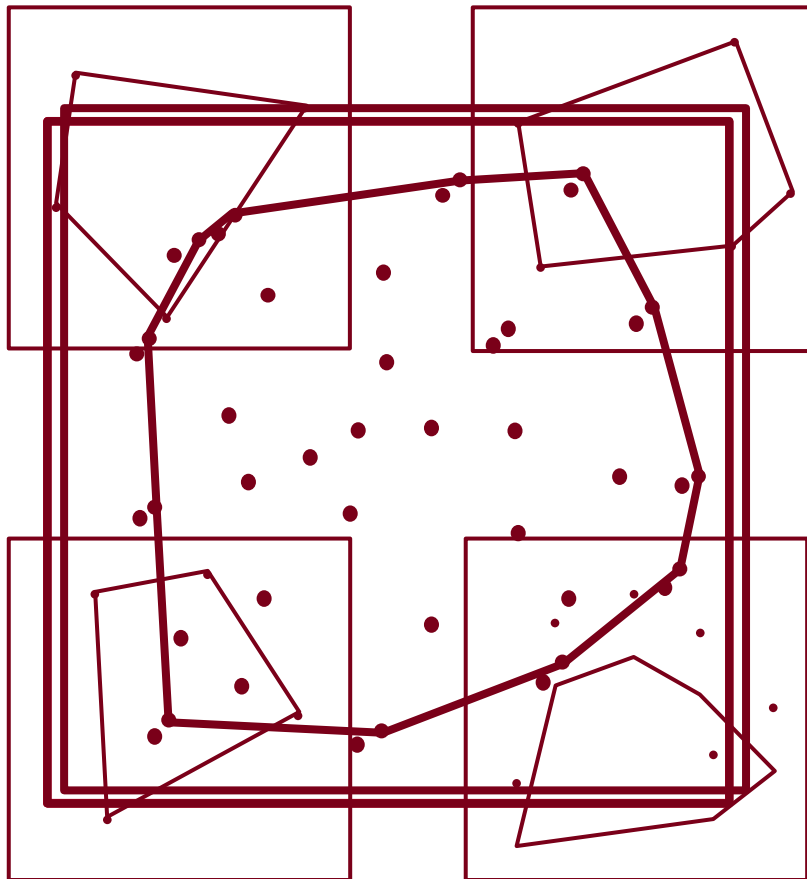


Convex Hull in CG_Hadoop

Hadoop

SpatialHadoop

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



Advanced Analytics

Partitioning

Local VD

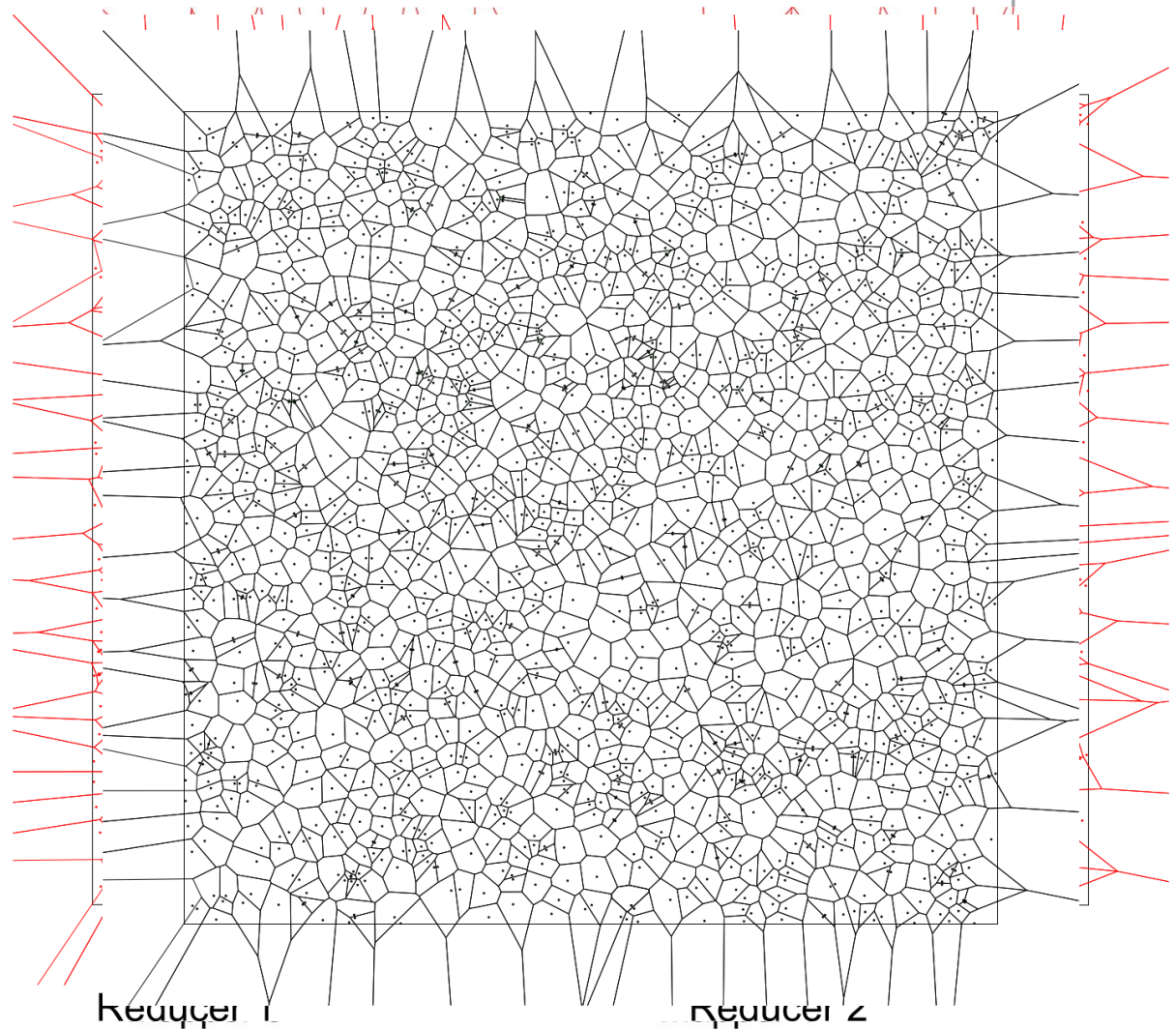
Pruning

Vertical Merge

Pruning

Horizontal Merge

Final output



Agenda

- › The ecosystem of SpatialHadoop
 - › Motivation
 - › Internal system design
 - › Applications
 - › Related work
 - › Performance Results
- › Interactive data exploration

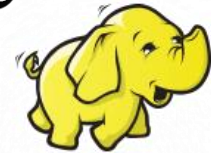
Other Big Spatial Data Systems

Parallel

SECONDO

MD-HBase

GeoSpark



Hadoop-GIS
Spatial Big Data Solutions



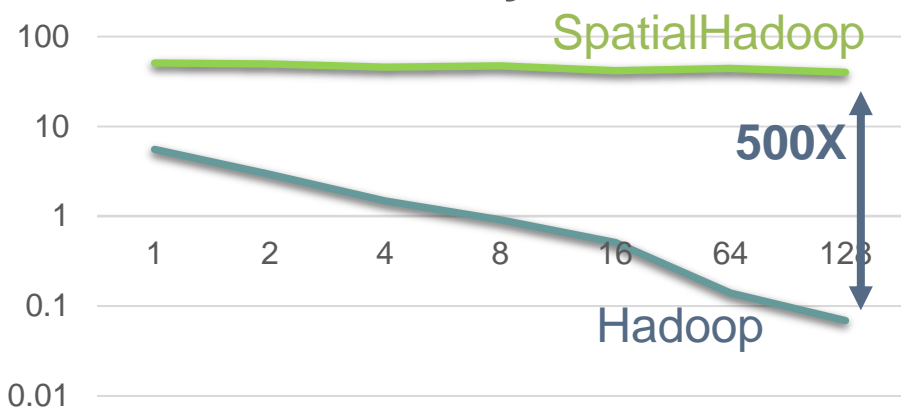
ESRI Tools
for Hadoop

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

Performance Results

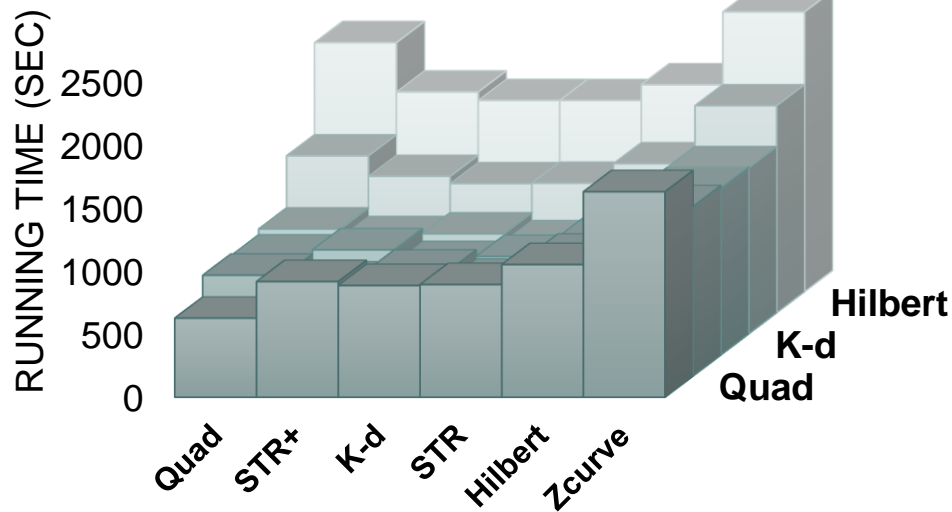


Throughput of Range Query

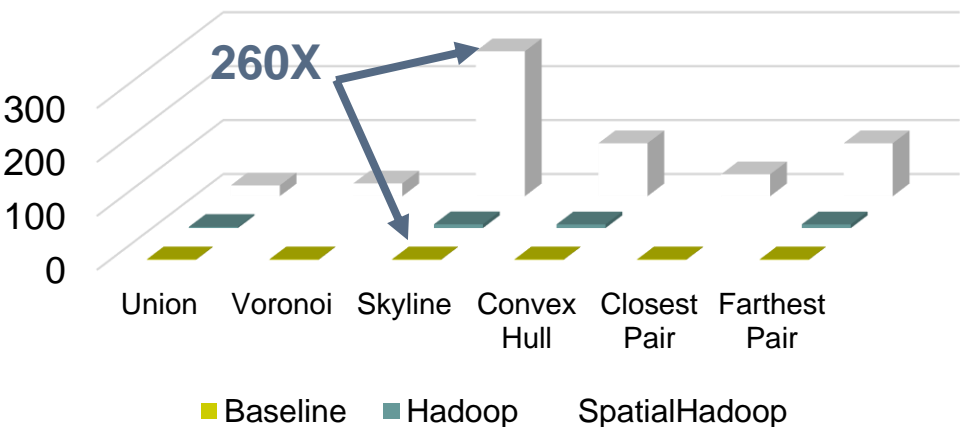


Spatial Join

Running time with different indexes



Speedup of CG_Hadoop



Agenda

- ~~The ecosystem of SpatialHadoop~~
 - ~~Motivation~~
 - ~~System design~~
 - ~~Applications~~
 - ~~Related work~~
 - ~~Performance results~~
- **Interactive data exploration**

Rise of Big Open Data

The home of the U.S. Government

Here you will find data, tools, and resources to connect and apply

BETA This is a new service – your [feedback](#) will help us to improve it



Find open data

World Bank Open Data

Free and open access to global development data

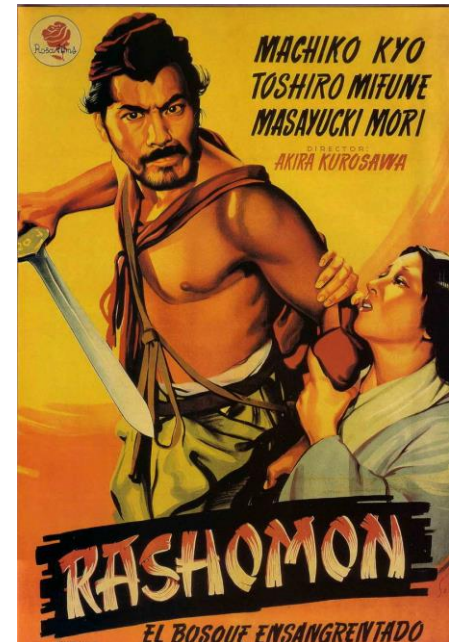
Search data e.g. GDP, population, Indonesia

Browse by [Country](#) or [Indicator](#)

most recent seven days. Data is extracted from the Chicago Police Department's (Citizen Law Enforcement Analysis and Reporting) system. In order

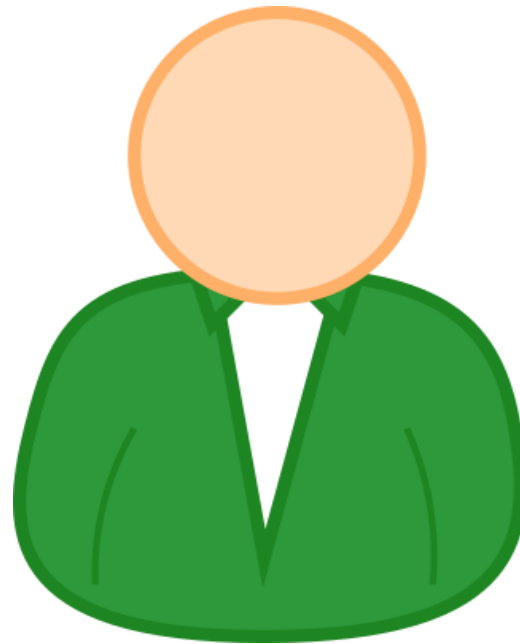
[More](#)

Did these data repositories work as expected?



The home of the U.S. Government's open data

Here you will find data, tools, and resources to conduct research, develop web and mobile applications, design data visualizations, and [more](#).

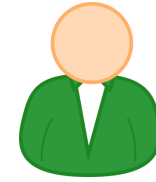


Helped the police department publish their data on Data.gov

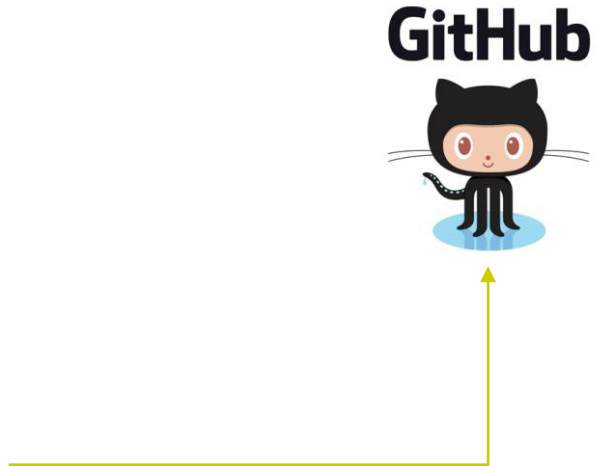
The Politician



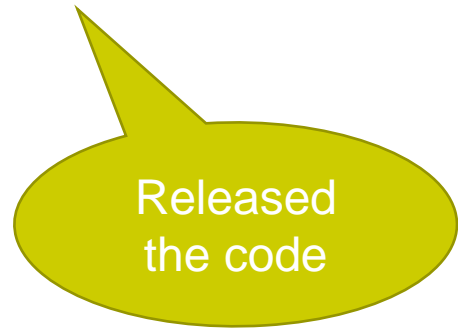
The Politician



Published a paper



GitHub



Released the code



Developed new spatial analytics algorithms to study event data

The Computer Scientist

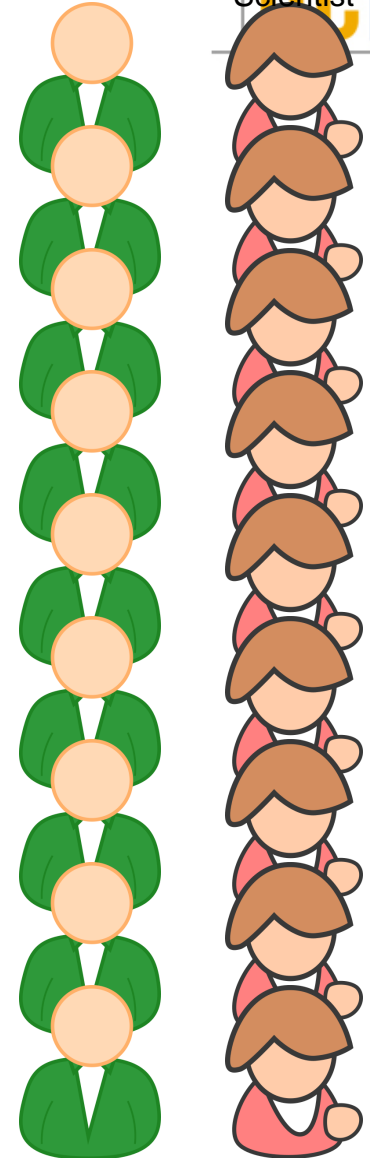
Want to study the relationship between crimes, demographics, and homelessness using data science

????

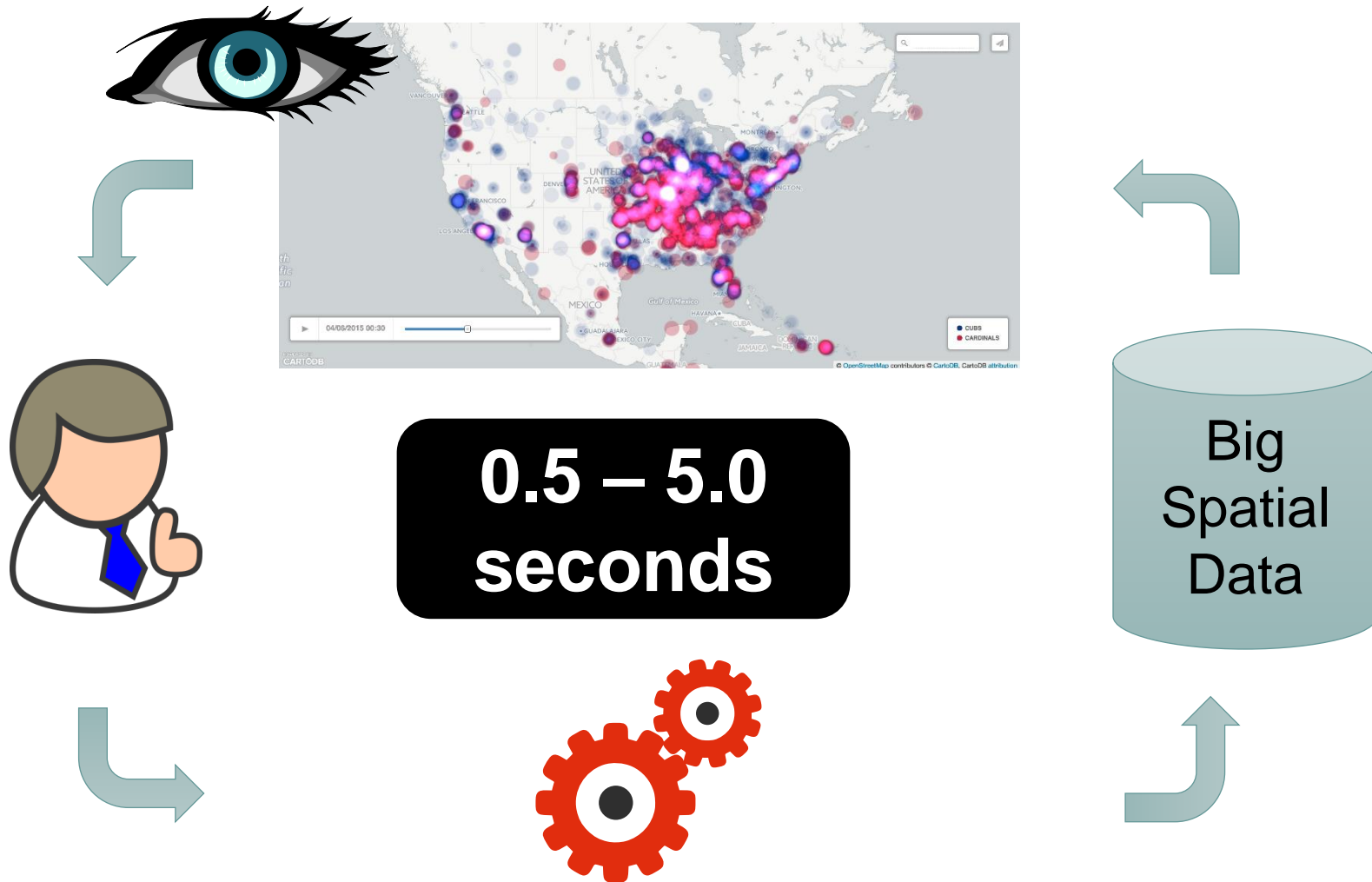
The Domain Scientist

The Politician

The Computer Scientist



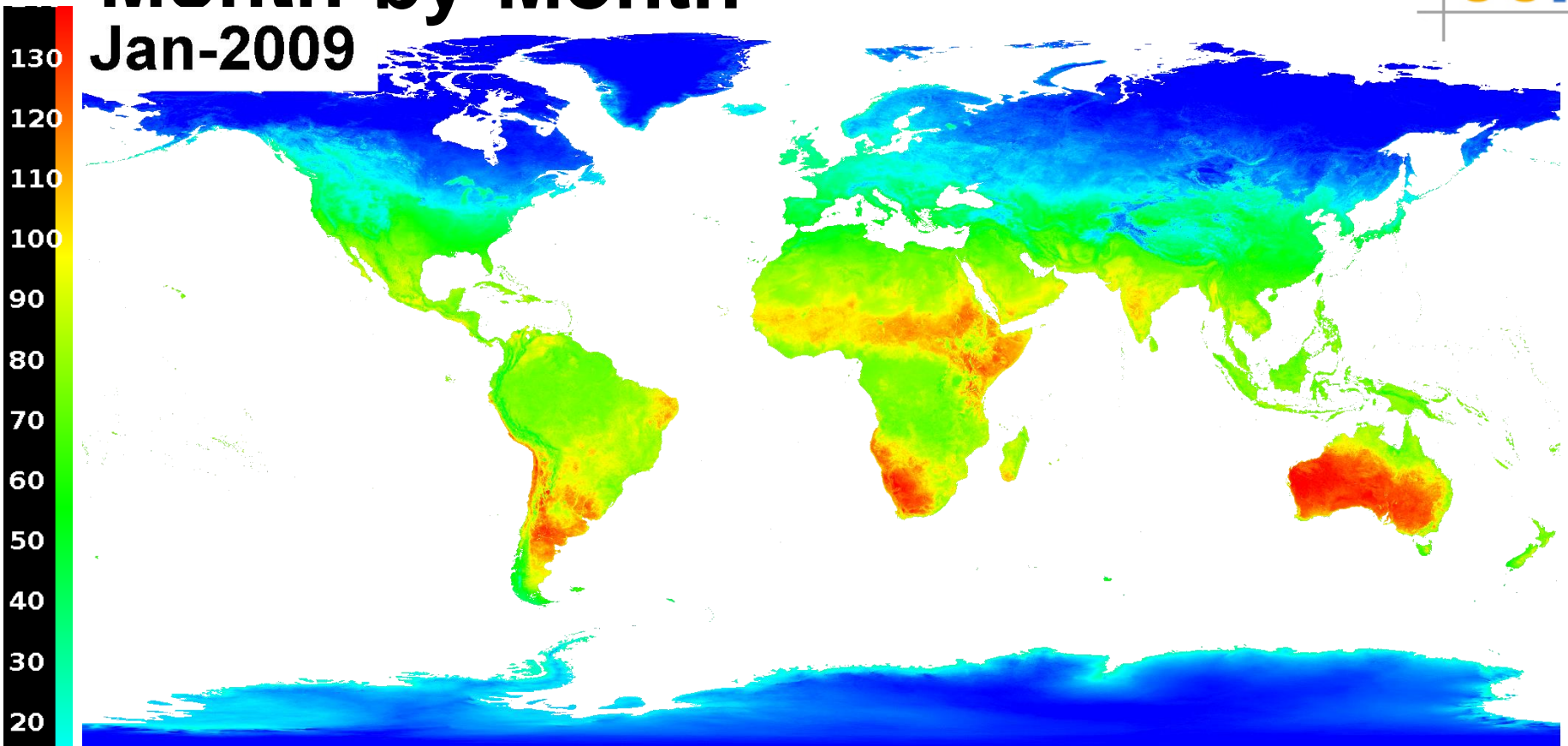
Exploratory Analytics



Heat Map From 2009 to 2014

Month-by-Month

Jan-2009



72 Frames × 14 Billion points per frame

Total = **1 Trillion points**

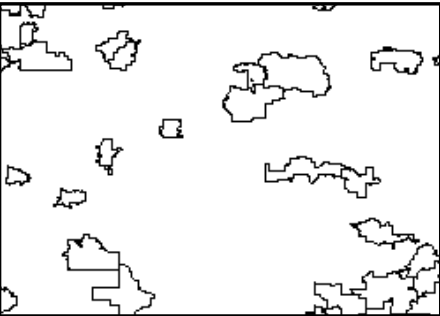
Created in **3 hours** on **10 nodes** instead of **60 hours**

Single Level Image

Input

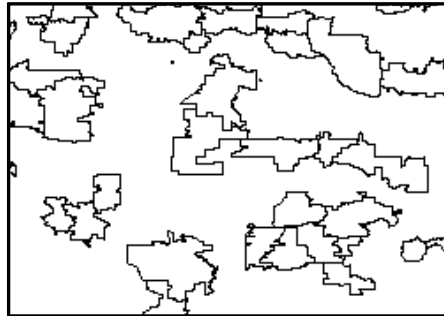
Split

caeaeeizester



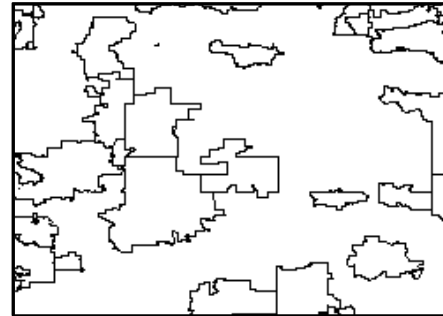
Split

caeaeeizester



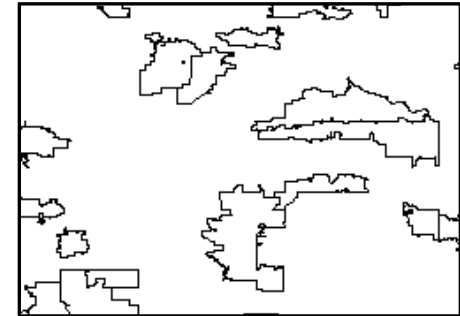
Split

caeaeeizester



Split

caeaeeizester



Merge
(Overlay)

Space Partitioning

Input

Split

Split

Split

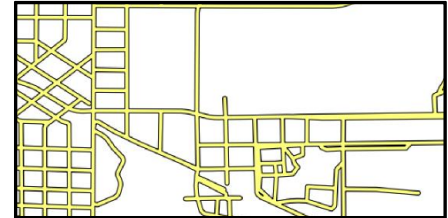
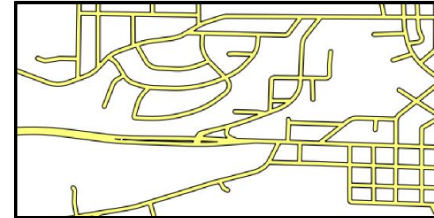
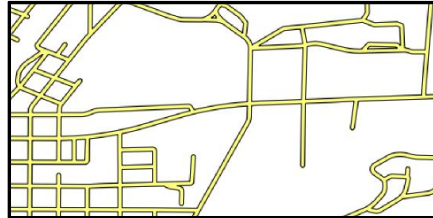
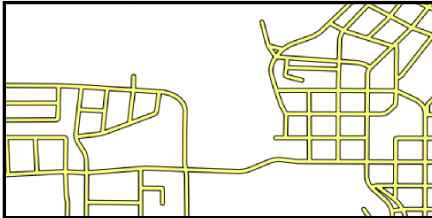
Split

`create-raster`

`create-raster`

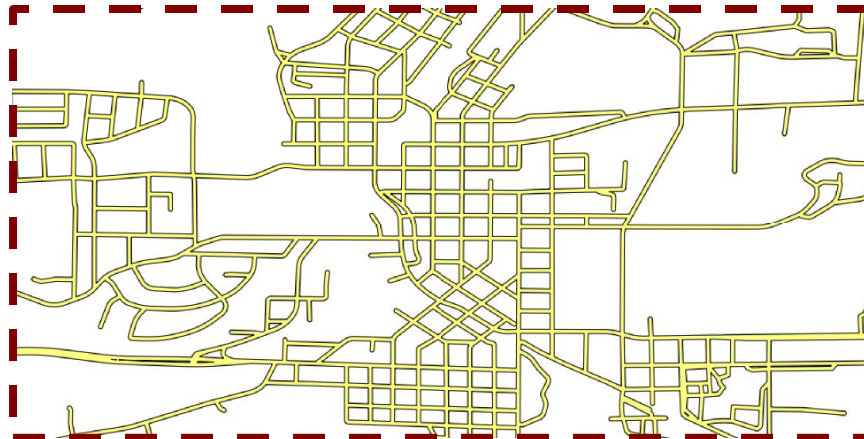
`create-raster`

`create-raster`

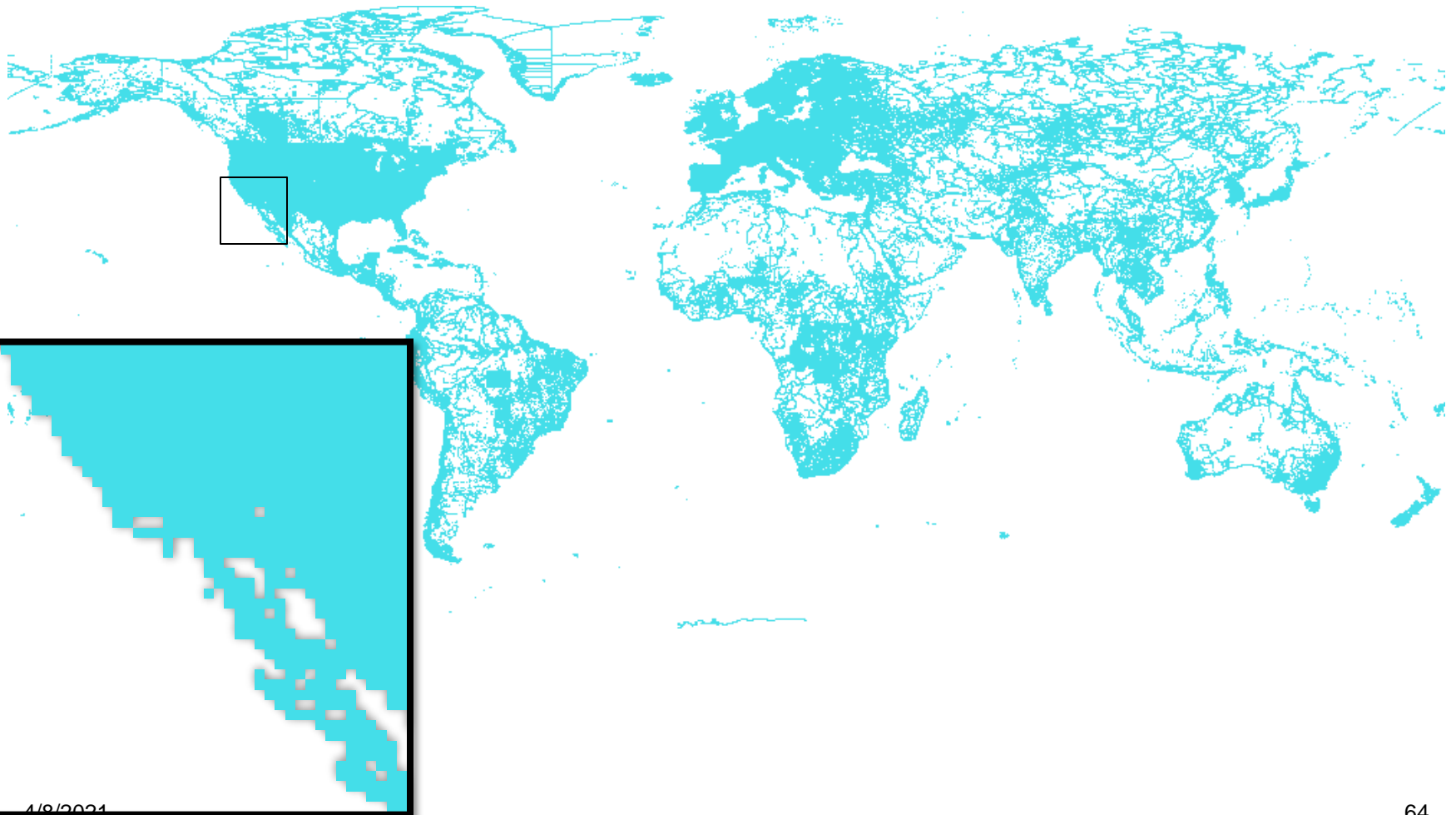


`create-raster`

Merge (Stitch)



Level of Details



Multilevel Visualization

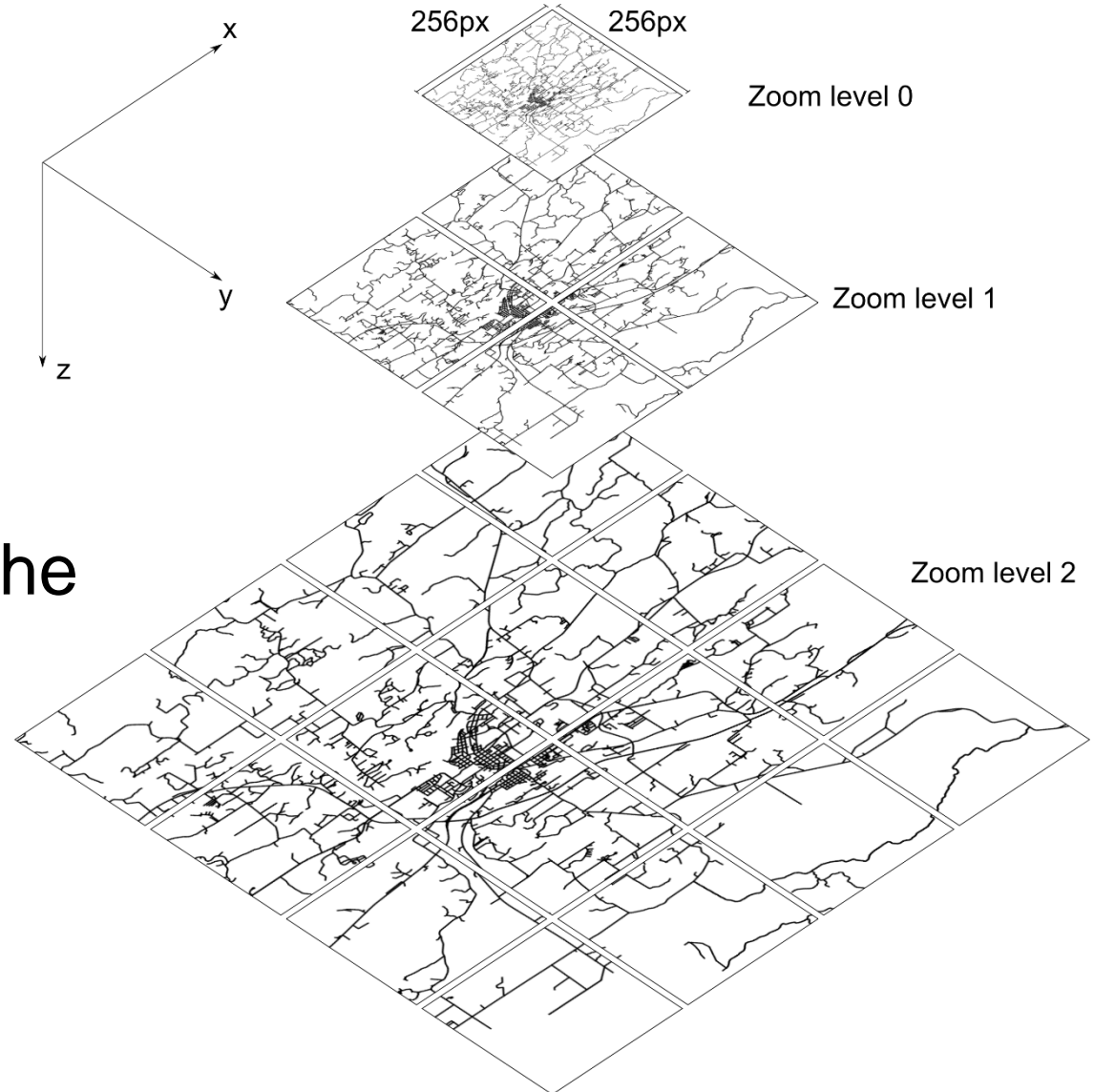


Map of California – 2GB

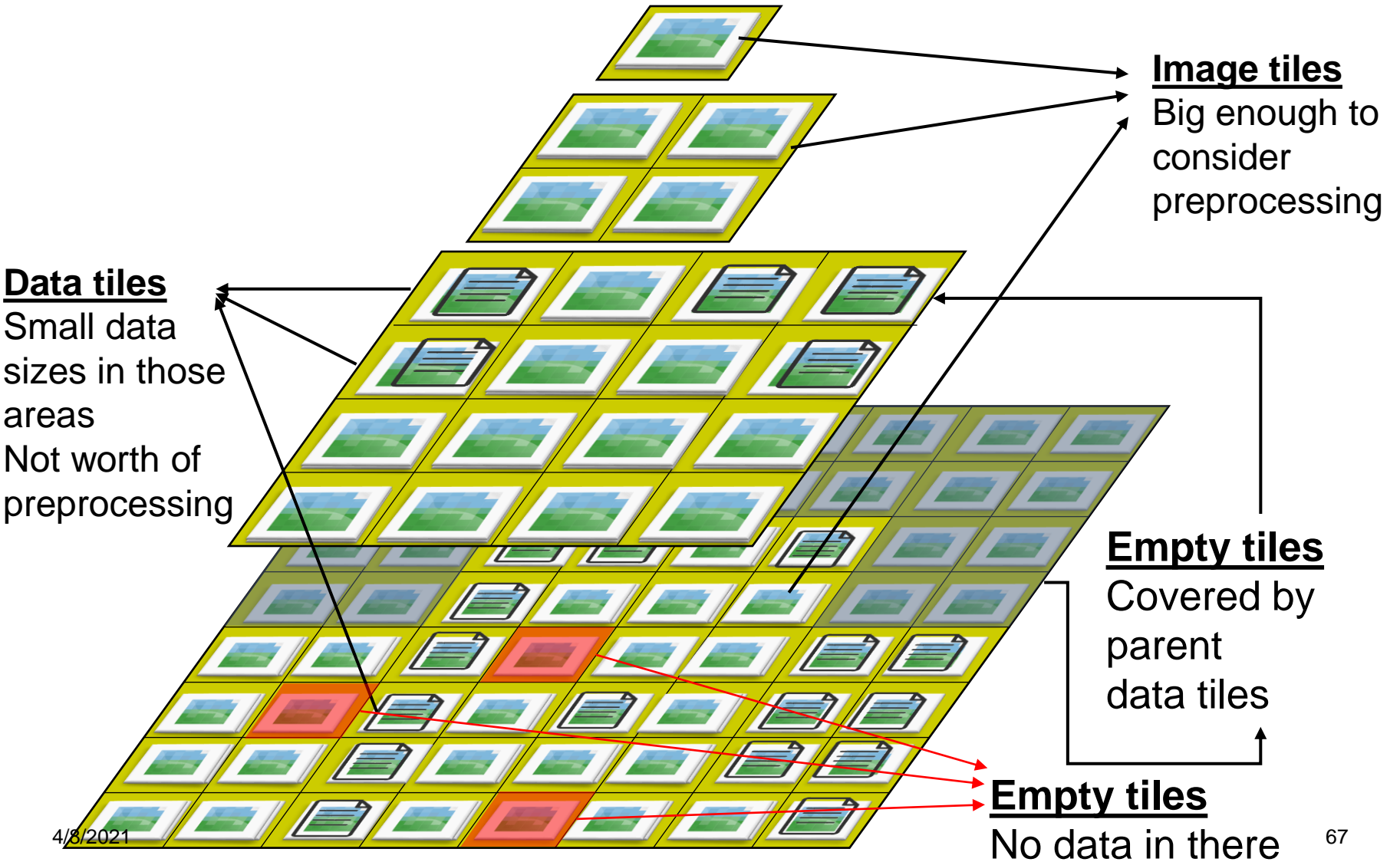
Generated in **2 minutes** on 10-node cluster using SpatialHadoop
instead of **one hour**

Multi-level Image

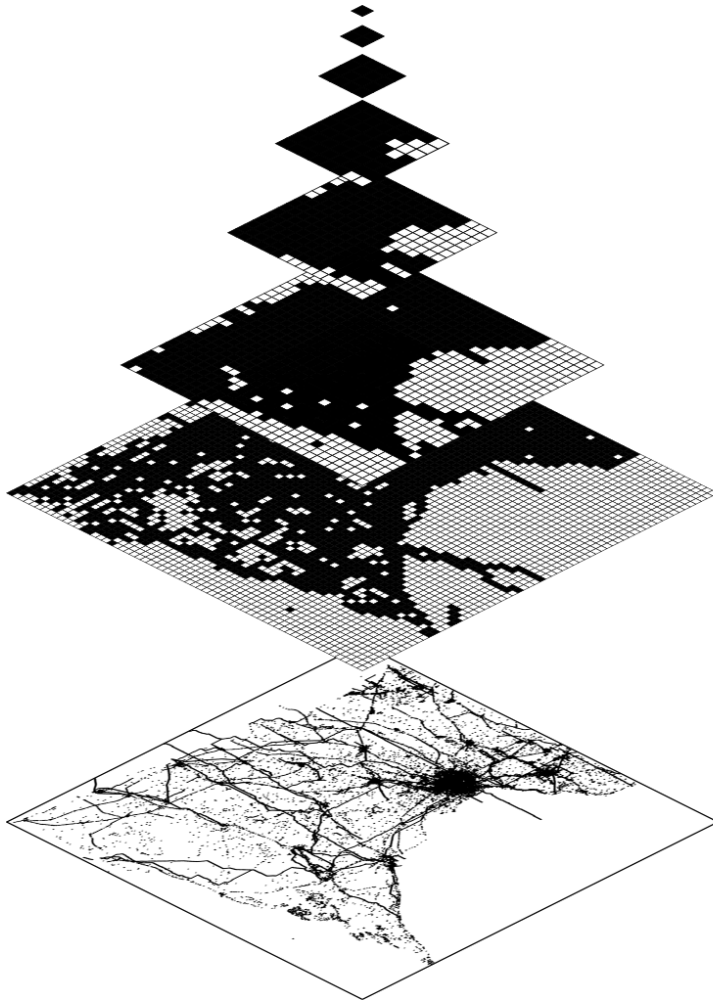
- ▶ Many images at different zoom levels
 - ▶ Pan
 - ▶ Zoom in/out
 - ▶ Fly to
- ▶ More details as the zoom level increases
- ▶ Number of tiles increases exponentially



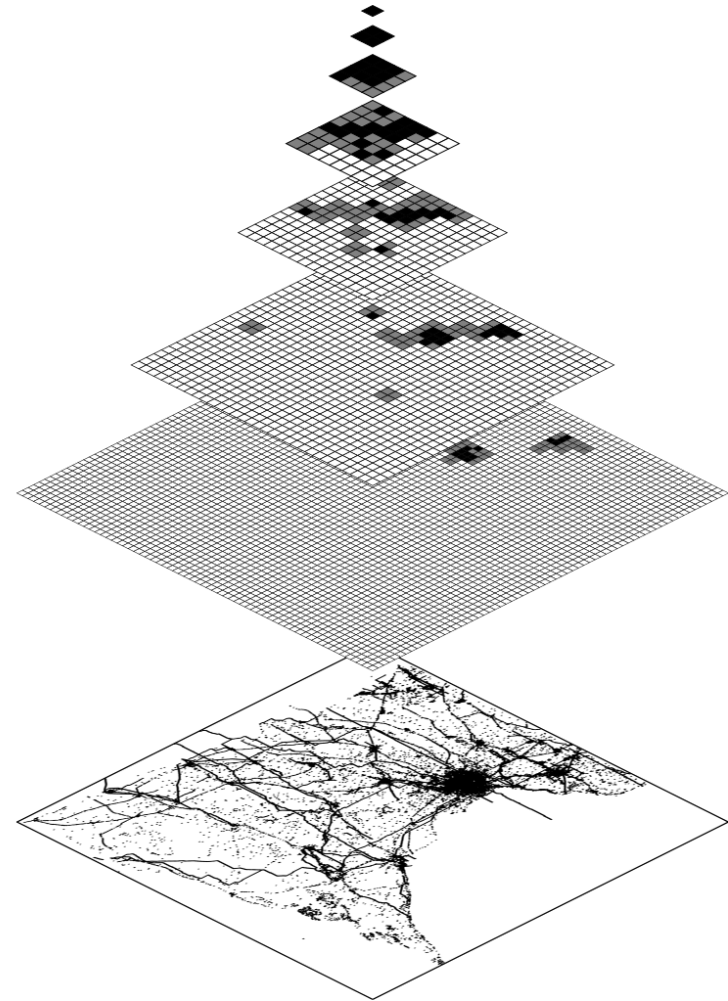
Adaptive Multilevel Visualization



Adaptive Multilevel Images



Full Image (3,160 tiles)



Adaptive Image (231 files)

Text search

232 Datasets

- MSBuildings
- Chicago Crimes
- eBird
- NE/countries
- NE/states_provinces
- NE/time_zones
- TIGER2018/CD
- TIGER2018/LINEARWATER
- TIGER2018/ROADS
- TIGER2018/ZCTA5

MSBuildings

[Project homepage](#) [Download data](#)

Size: 20.6 GB

Number of records: 125 m

Number of points: 753 m

Format:

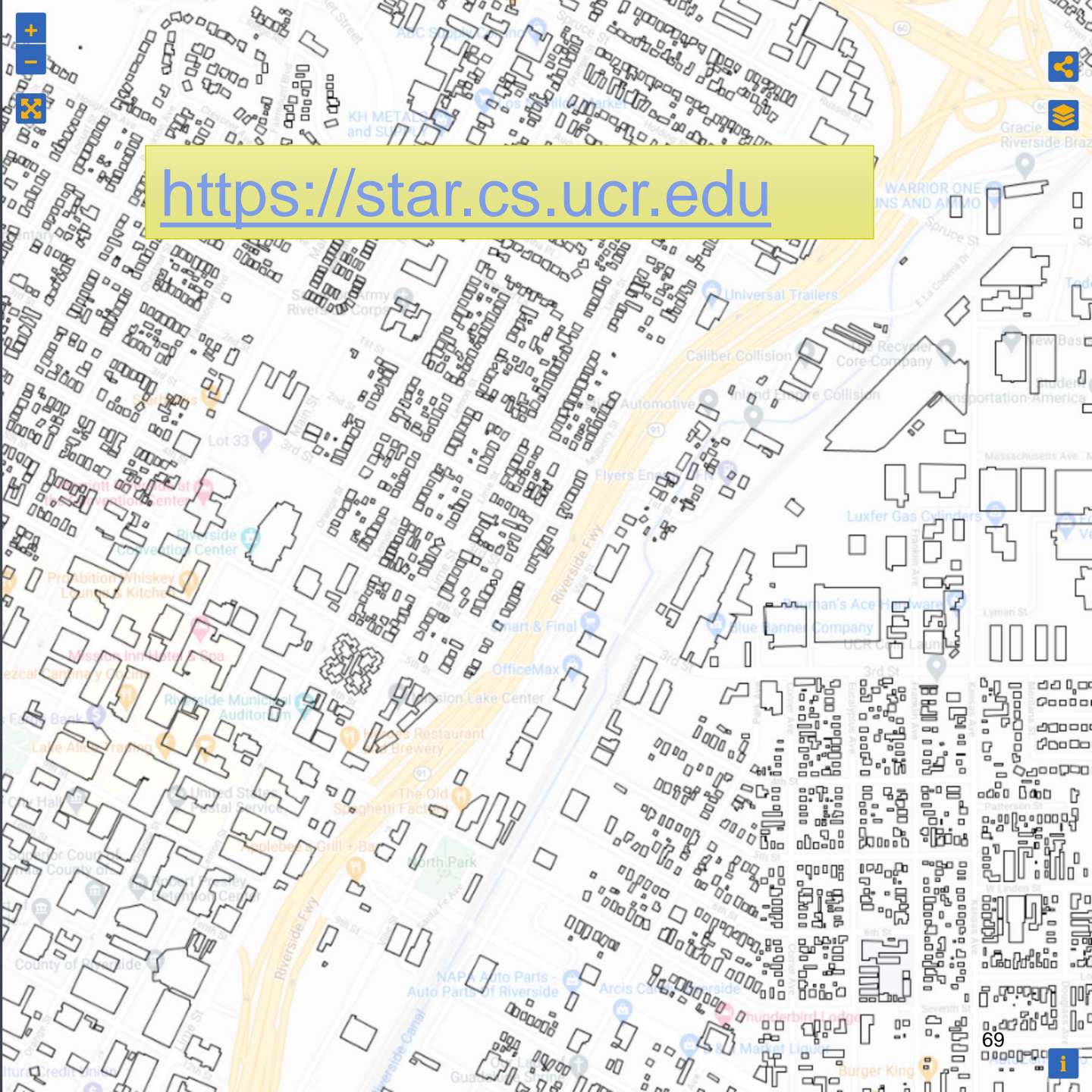
Geometry type: POLYGON

Description: This dataset contains 125,192,184 computer generated building footprint in all 50 US states. Provided by Microsoft.

Attributes:

geometry

<https://star.cs.ucr.edu>



Summary

Apps

Visualization

Operations

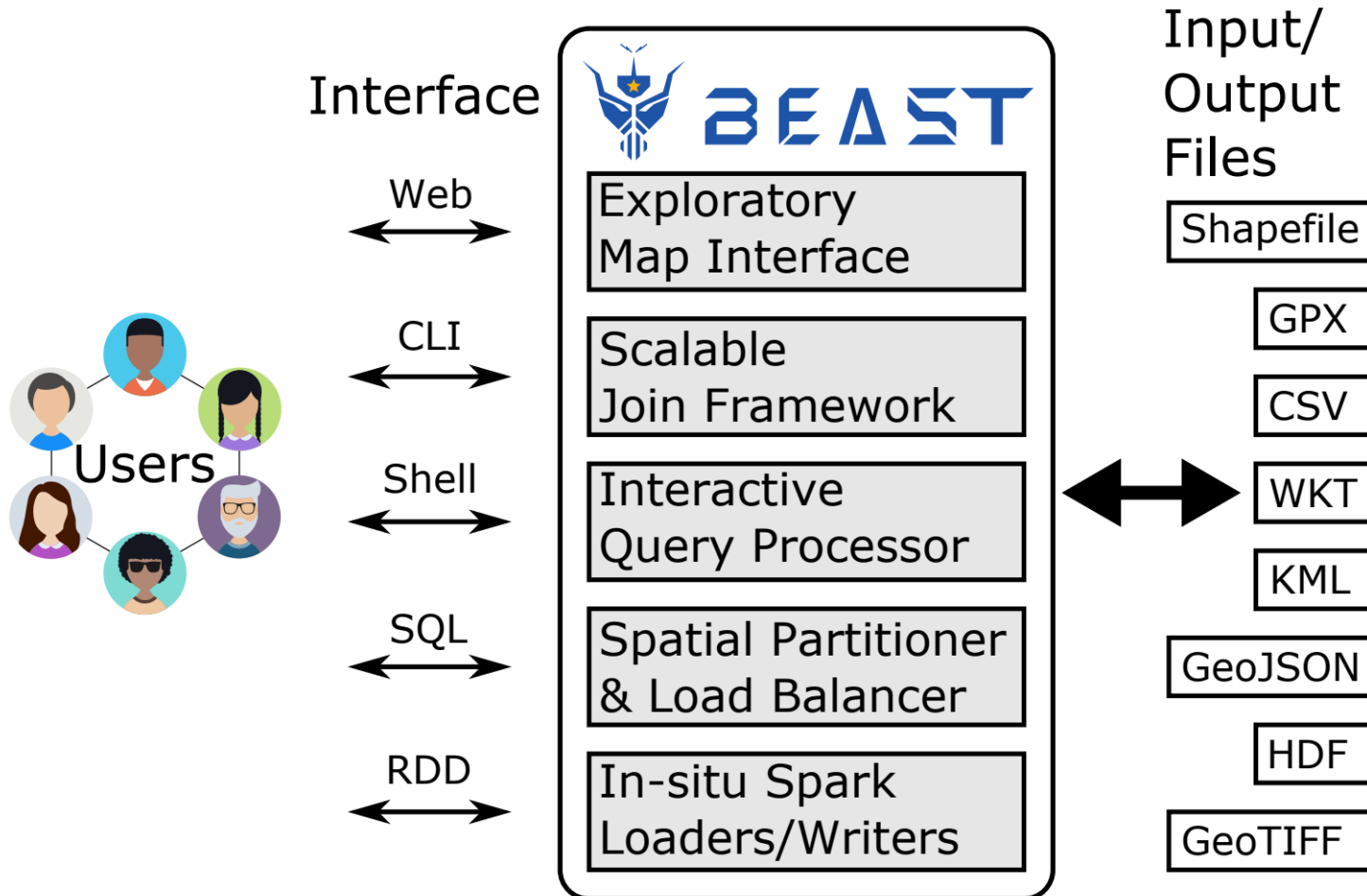
Indexes

Future Work

- More interaction in visualization
 - Non-spatial filters, e.g., temporal
 - Select individual objects on the map (on big data of course)
- Approximate query processing on big spatial data
- Utilize machine learning for query optimization and better user experience
- Support incremental updates to datasets

Beast

› Big Exploratory Analytics on Spatio-Temporal Data



Thank You

Questions?