

# State of Geospatial BigData

Mansour Raad  
thunderheadxplorer.blogspot.com  
mraad@esri.com  
@mraad



“WHERE” IS UBIQUITOUS !

- **Where** is the closest ATM ?
- **Where** is the best location to place my store ?
- **Where** is UBL ?
- **Where** is the next Ebola/Zika outbreak ?



# A BIT OF HISTORY...

With Esri Specifically :-)



~ 1990

# Shapefile

From Wikipedia, the free encyclopedia

The **Esri shapefile**, or simply a **shapefile**, is a popular geospatial vector [data format for geog](#) developed and regulated by [Esri](#) as a (mostly) [open specification](#) for data interoperability amo Shapefiles spatially describe [vector](#) features: [points](#), [lines](#), and [polygons](#), representing, for exa usually has [attributes](#) that describe it, such as *name* or *temperature*.

## Contents [\[hide\]](#)

### 1 Overview

- 1.1 Shapefile shape format (.shp)
- 1.2 Shapefile shape index format (.shx)
- 1.3 Shapefile attribute format (.dbf)
- 1.4 Shapefile spatial index format (.sbn)

### 2 Limitations

- 2.1 Topology and shapefiles
- 2.2 Spatial representation
- 2.3 Data storage
- 2.4 Mixing shape types

### 3 See also

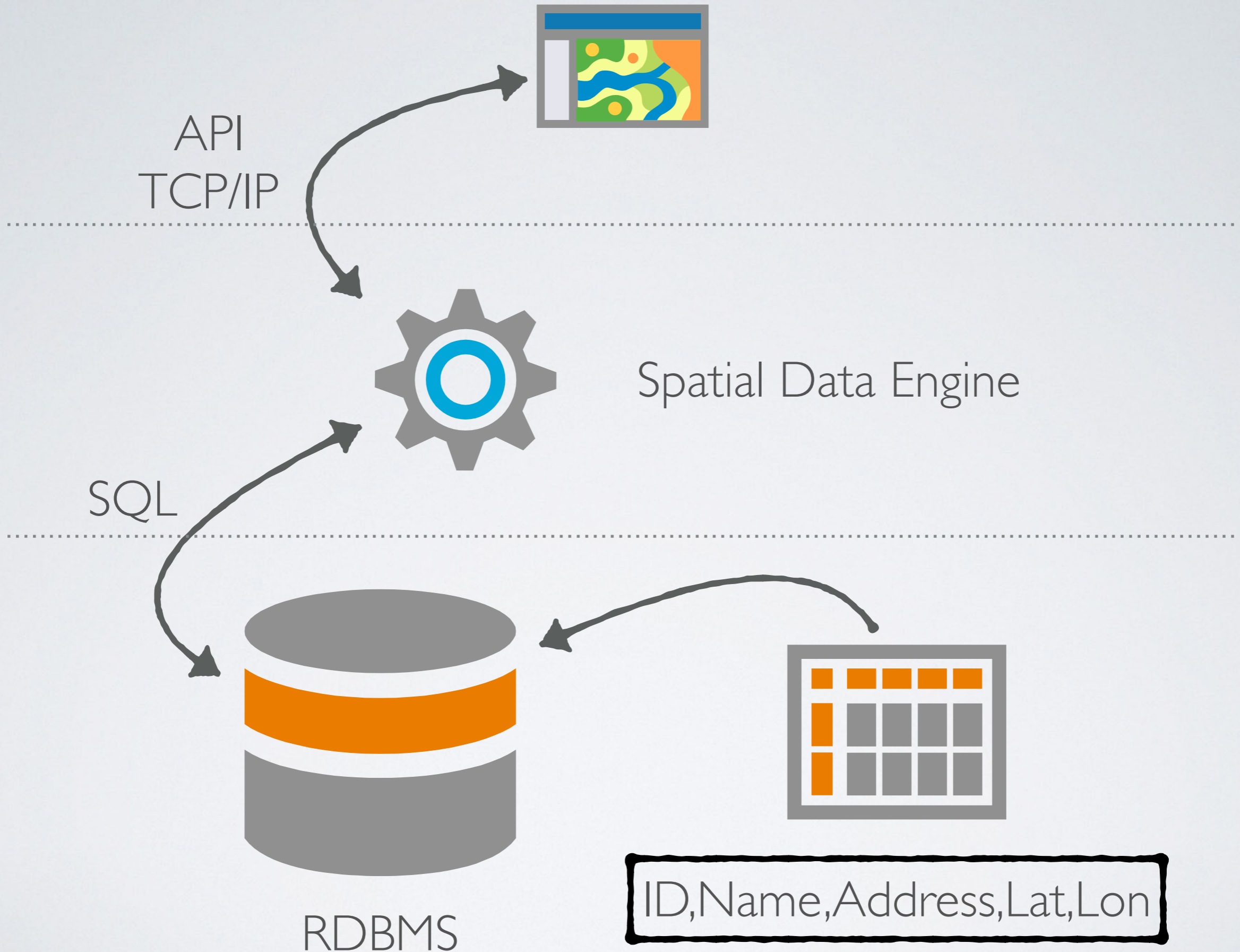
### 4 References

### 5 External links

<http://en.wikipedia.org/wiki/Shapefile>

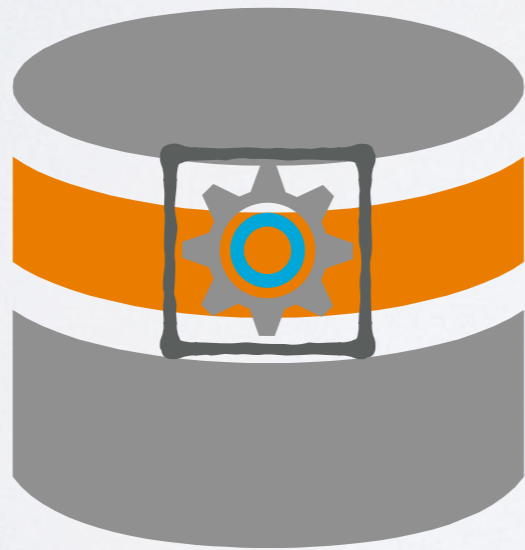
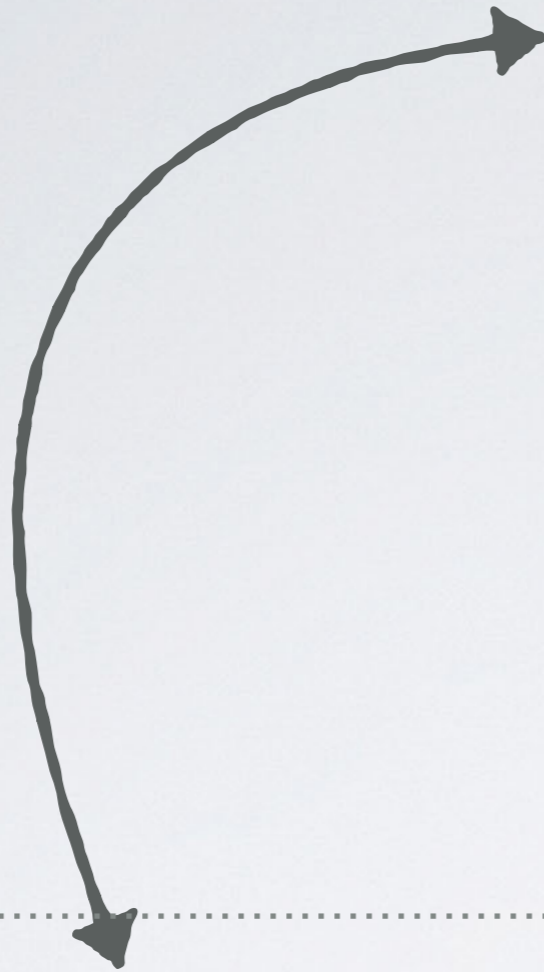
1995



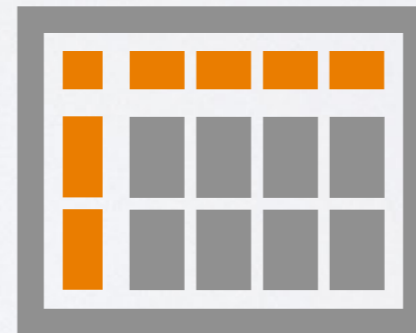


> 1996-2005

xDBC



RDBMS



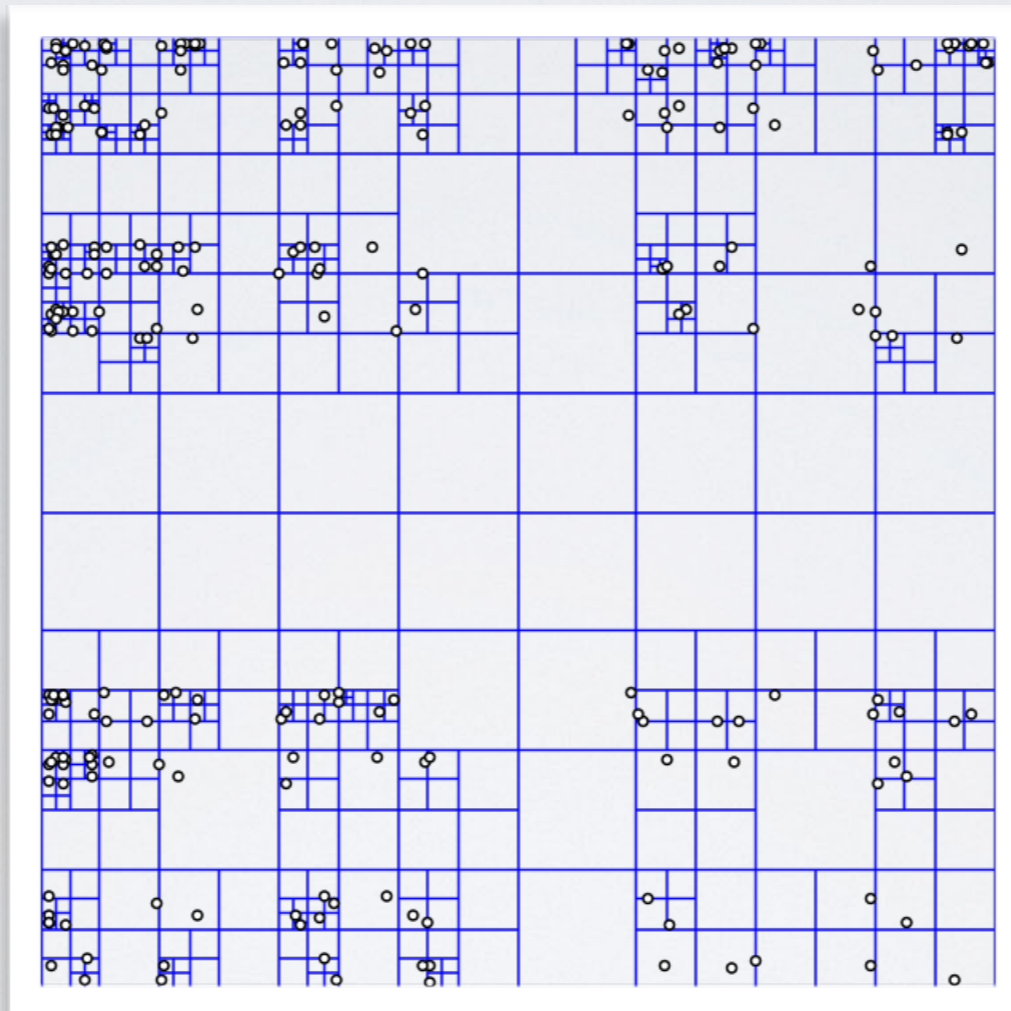
ID,Name,Address,Lat,Lon



# SPATIAL INDEXING

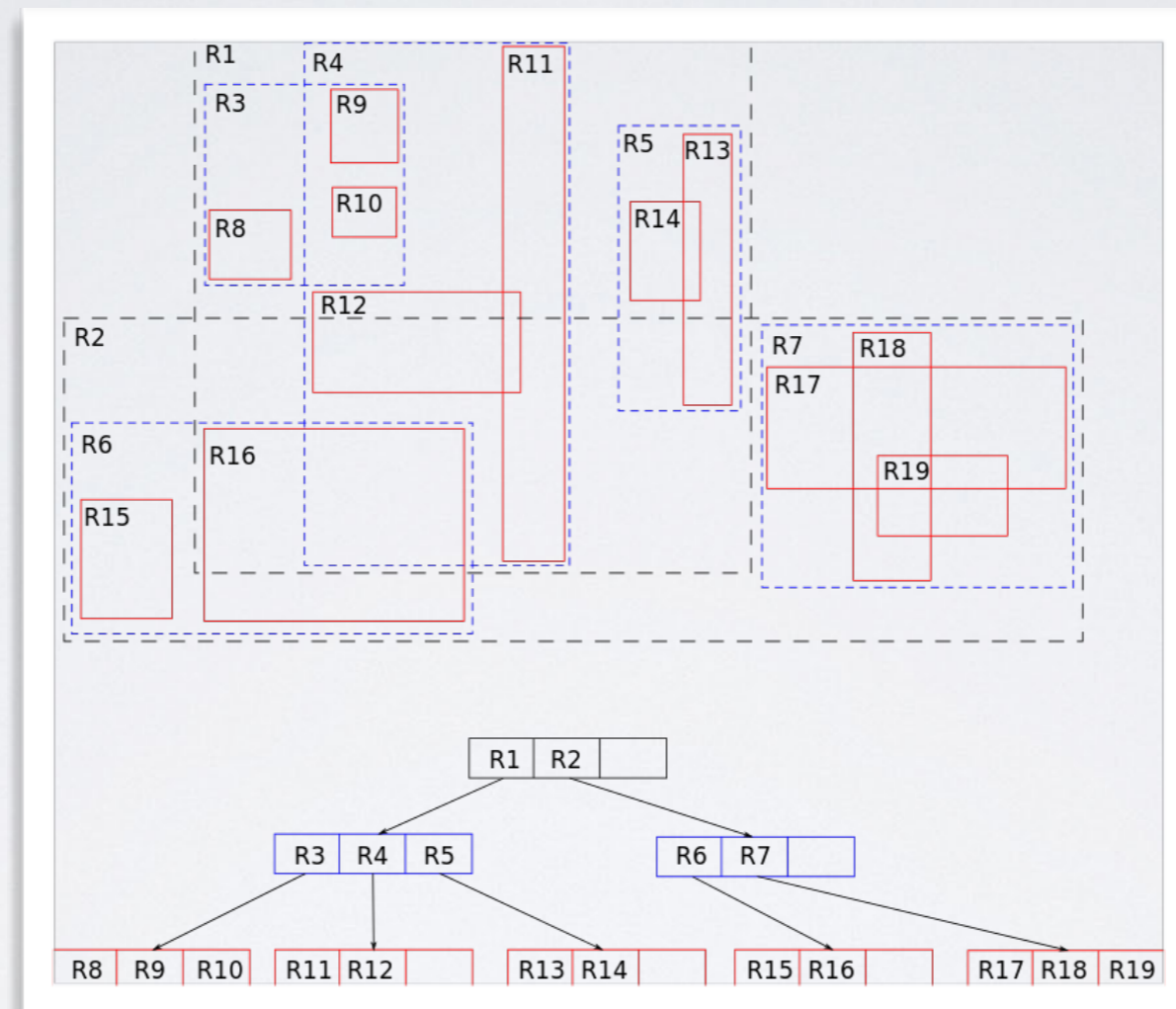
[http://en.wikipedia.org/wiki/Spatial\\_database#Spatial\\_index](http://en.wikipedia.org/wiki/Spatial_database#Spatial_index)

# QUADTREE



<http://en.wikipedia.org/wiki/Quadtree>

# R-TREE



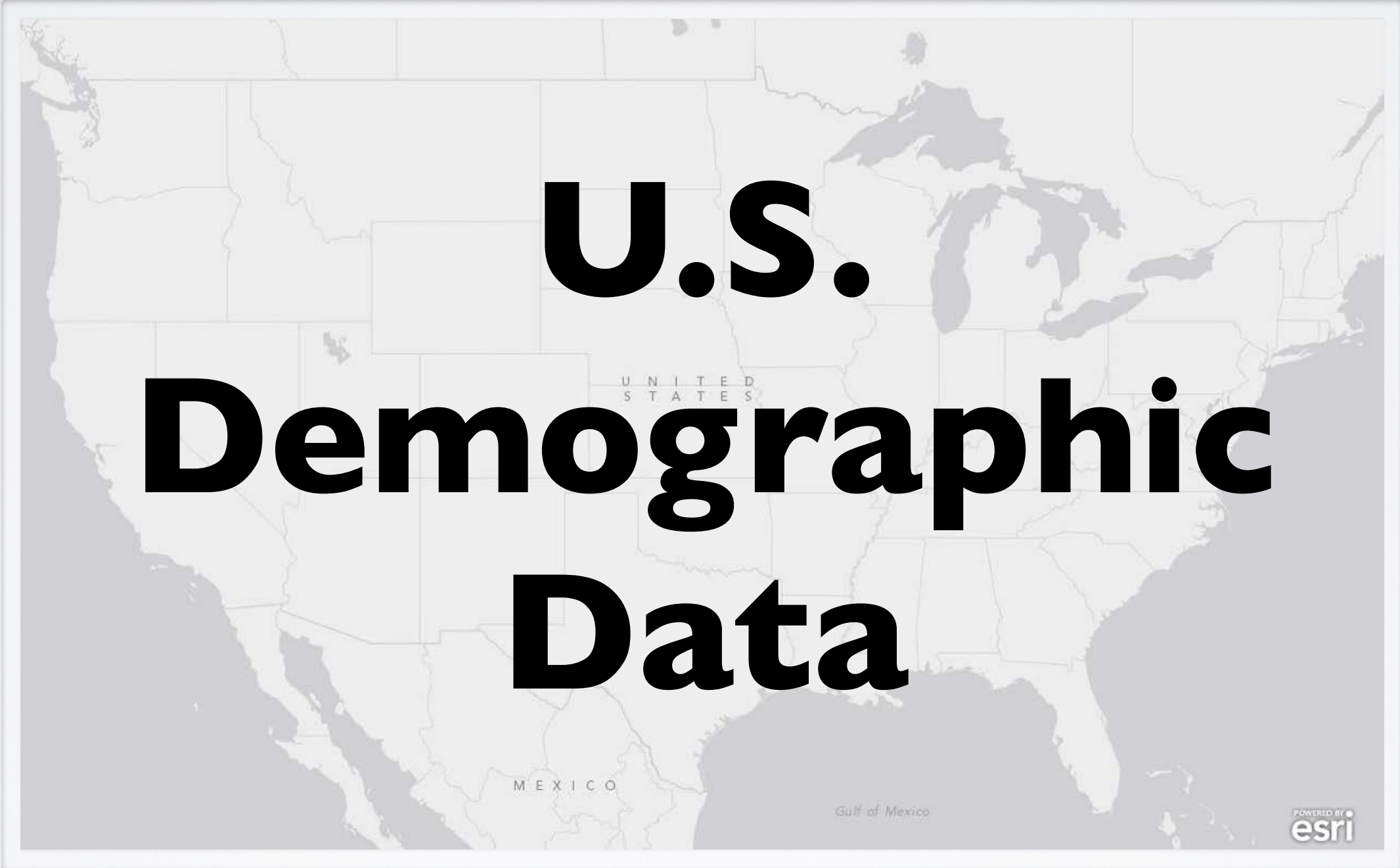
<http://en.wikipedia.org/wiki/R-tree>



(NOT SO) MODERN DAY...

STORY TIME...





# U.S. Demographic Data

POWERED BY  
esri





**FOR EACH LOCATION  
FOR EACH DEMOGRAPHIC**



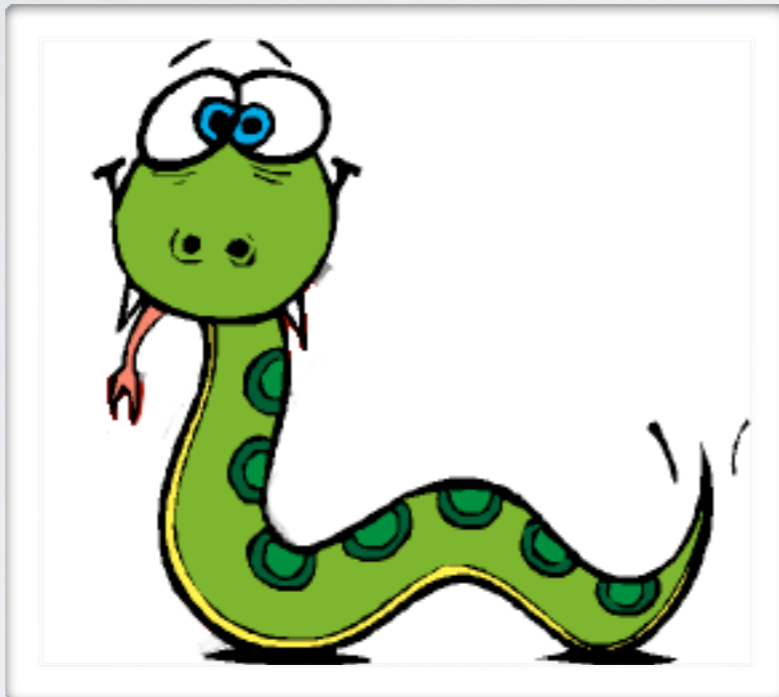
**50 MILE HEATMAP**







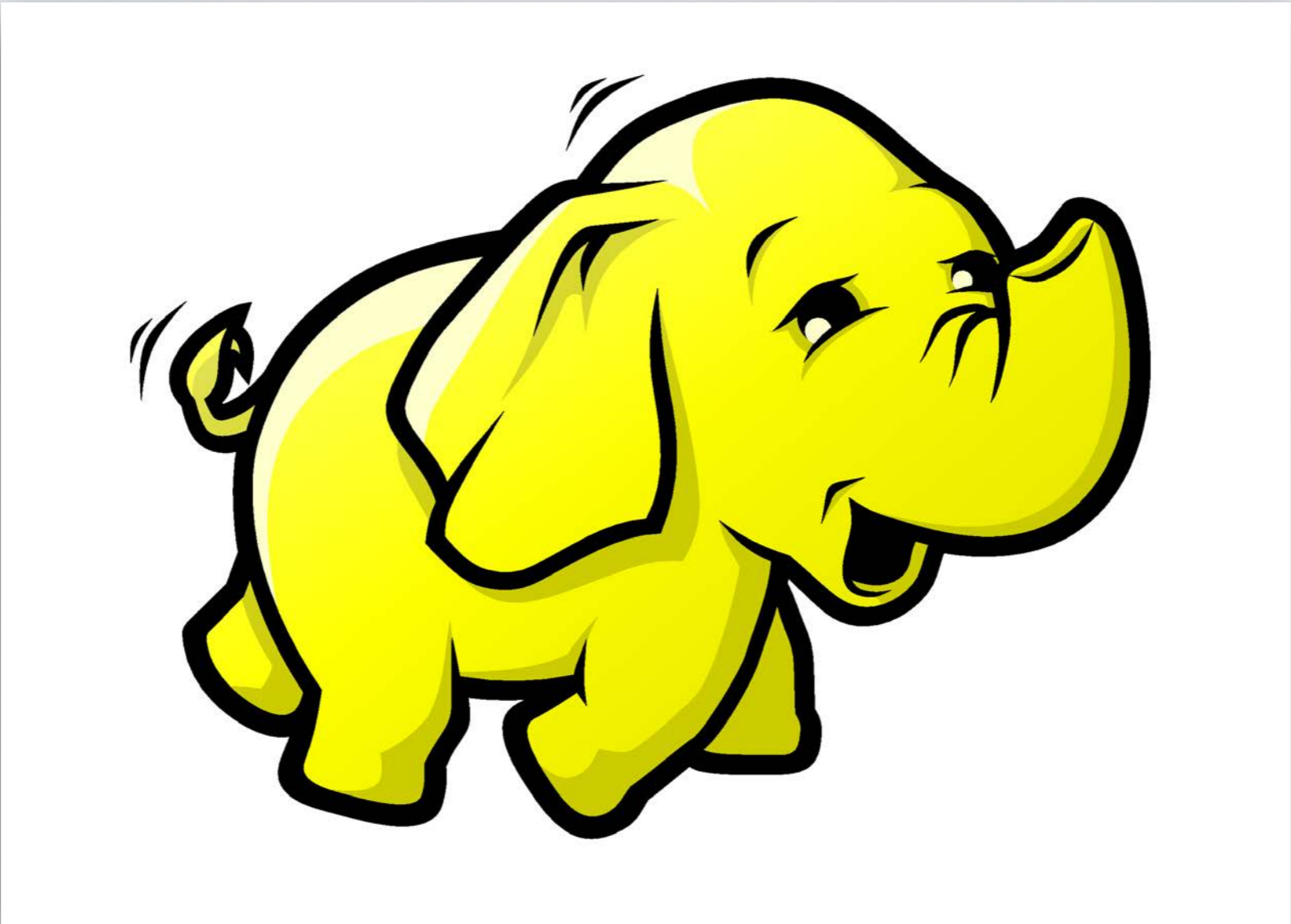
# TRADITIONAL MEANS...



**14 Days**

**850 GB Raster Files**

**BETTER WAY ?**







**amazon**  
**web services™**

2012-07-18 15:41 EDT

0 hours 2 minutes

21

### Create a New Job Flow Cancel X

DEFINE JOB FLOW   SPECIFY PARAMETERS   CONFIGURE EC2 INSTANCES   ADVANCED OPTIONS   BOOTSTRAP ACTIONS   REVIEW

Please review the details of your job flow and click "Create Job Flow" when you are ready to launch your Hadoop Cluster.

<b>Job Flow Name:</b>	My Heat Map	
<b>Type:</b>	Custom Jar	<a href="#">Edit Job Flow Definition</a>
<b>Jar Location:</b>	s3n://ags/MRHeatMap.jar	
<b>Jar Arguments:</b>	<input type="text" value="s3n://mraads3 s3n://mraadout/heatmap"/>	<a href="#">Edit Job Flow Parameters</a>
<b>Master Instance Type:</b>	m1.small	
<b>Core Instance Type:</b>	m1.large	<b>Instance Count: 10</b> <a href="#">Edit EC2 Configs</a>
<b>Amazon EC2 Key Pair:</b>		
<b>Amazon Subnet Id:</b>		
<b>Amazon S3 Log Path:</b>	s3n://mraadlogs/mrheatmaplog	
<b>Enable Debugging:</b>	No	
<b>Keep Alive:</b>	No	
<b>Termination Protected:</b>	No	<a href="#">Edit Advanced Options</a>
<b>Bootstrap Actions:</b>	No Bootstrap Actions created for this Job Flow	<a href="#">Edit Bootstrap Actions</a>

[Back](#) [Create Job Flow](#)

**Note:** Once you click "Create Job Flow," instances will be launched and you will be charged accordingly.



services

edit shortcut

### Your Elastic MapReduce Job Flows

Region:



US East (Virginia)



Create New Job Flow



Terminate



Debug

Viewing:

All

	Name	State	Creation Date	Elapsed Time	Normalized Instance Hours
<input type="checkbox"/>	HeatMap4	COMPLETED	2012-07-19 07:19 EST	0 hours 30 minutes	41
<input type="checkbox"/>	HeatMap31	COMPLETED	2012-07-18 16:07 EST	0 hours 5 minutes	21



Outstanding balance for this statement

\$37.67

## Details

[Expand All Services](#) | [Collapse All Services](#)

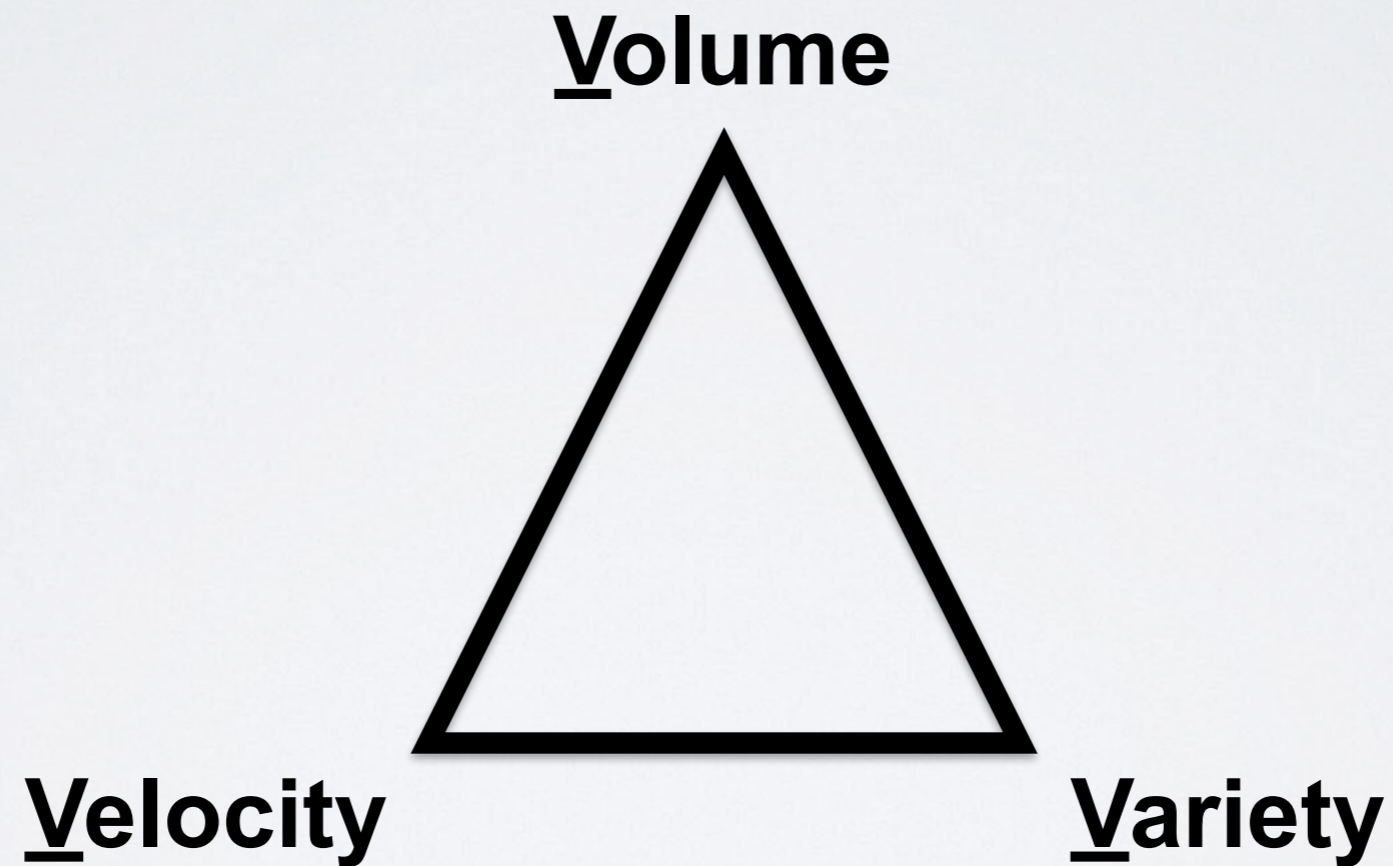
[Printer Friendly Version](#)

AWS Service Charges		\$37.67
<input type="checkbox"/> <b>Amazon Elastic Compute Cloud</b>		\$31.59
<a href="#">Download Usage Report »</a>		
<input type="checkbox"/> <b>Amazon SimpleDB</b>		\$0.00
<a href="#">Download Usage Report »</a>		
<input type="checkbox"/> <b>Amazon Simple Notification Service</b>		\$0.00
<a href="#">Download Usage Report »</a>		
<input type="checkbox"/> <b>Amazon Simple Storage Service</b>		\$2.00
<a href="#">Download Usage Report »</a>		
<input type="checkbox"/> <b>Amazon Elastic MapReduce</b>		\$4.07
<a href="#">Download Usage Report »</a>		
<input type="checkbox"/> <b>AWS Data Transfer (excluding Amazon CloudFront)</b>		\$0.01
<input type="checkbox"/> <b>VAT to be collected</b>		\$0.00

† Usage and recurring charges for this statement period will be charged on your next billing date, August 1, 2012. Estimated charges shown on this page, or shown on any notifications that we send to you, may differ from your actual charges for this statement period. This is because estimated charges presented on this page do not include usage charges accrued during this statement period after the date you view this page. Similarly, information about estimated charges sent to you in a notification do not include usage charges accrued during this statement period after the date you view the notification. One-time

**BIG DATA ?**

# U R IN BIGDATA SPACE IF...

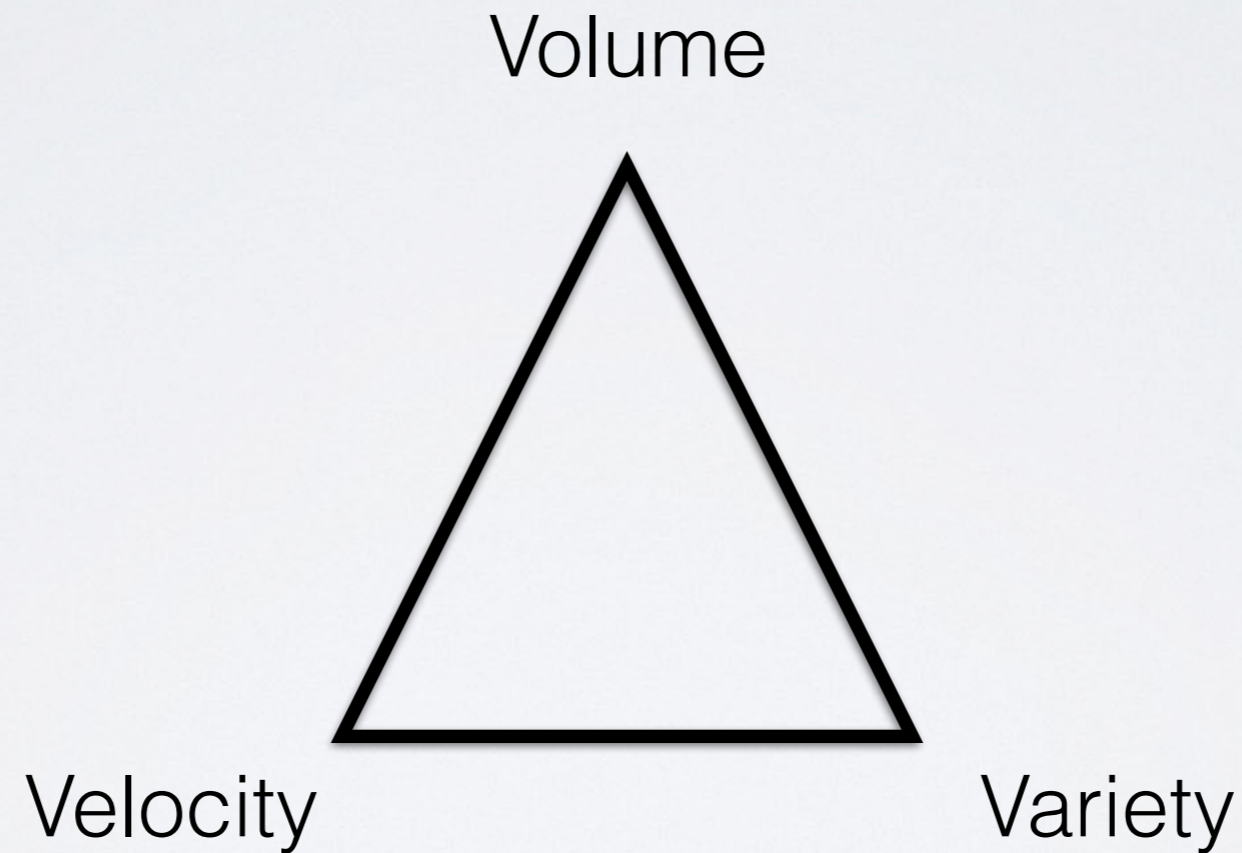




# BUT THEN I'VE SEEN...

- Volume** → data at rest
- Velocity** → data in motion
- Variety** → many types
- Veracity** → data in doubt
- Validity** → data that is correct
- Visualization** → data in patterns
- Vulnerability** → data at risk
- Value** → data that is meaningful

# I'M STICKING WITH...







NOSQL  
(NOT ONLY SQL :-)

# GEOJSON

<http://geojson.org/>

```
{
  "type": "Feature",
  "geometry": {
    "type": "Point",
    "coordinates": [125.6, 10.1]
  },
  "properties": {
    "name": "Dinagat Islands"
  }
}
```



- Points
- Lines
- Polygons
- Multipoints
- Multilines
- Multipolygons
- Geometry Collection

byte[]

byte[]

Sorted



Key1 → Value1

Key2 → Value2

...

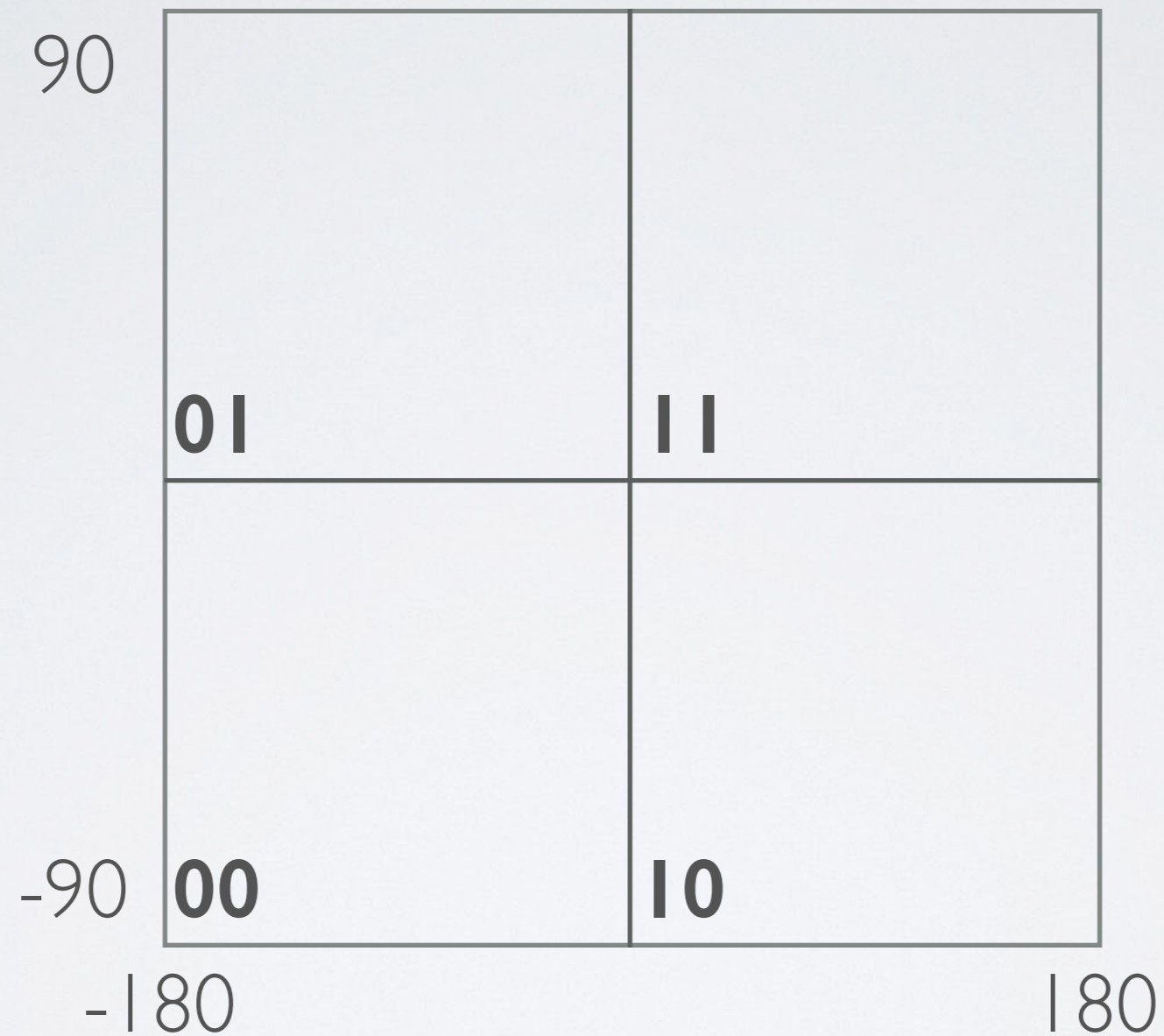
KeyN → ValueN

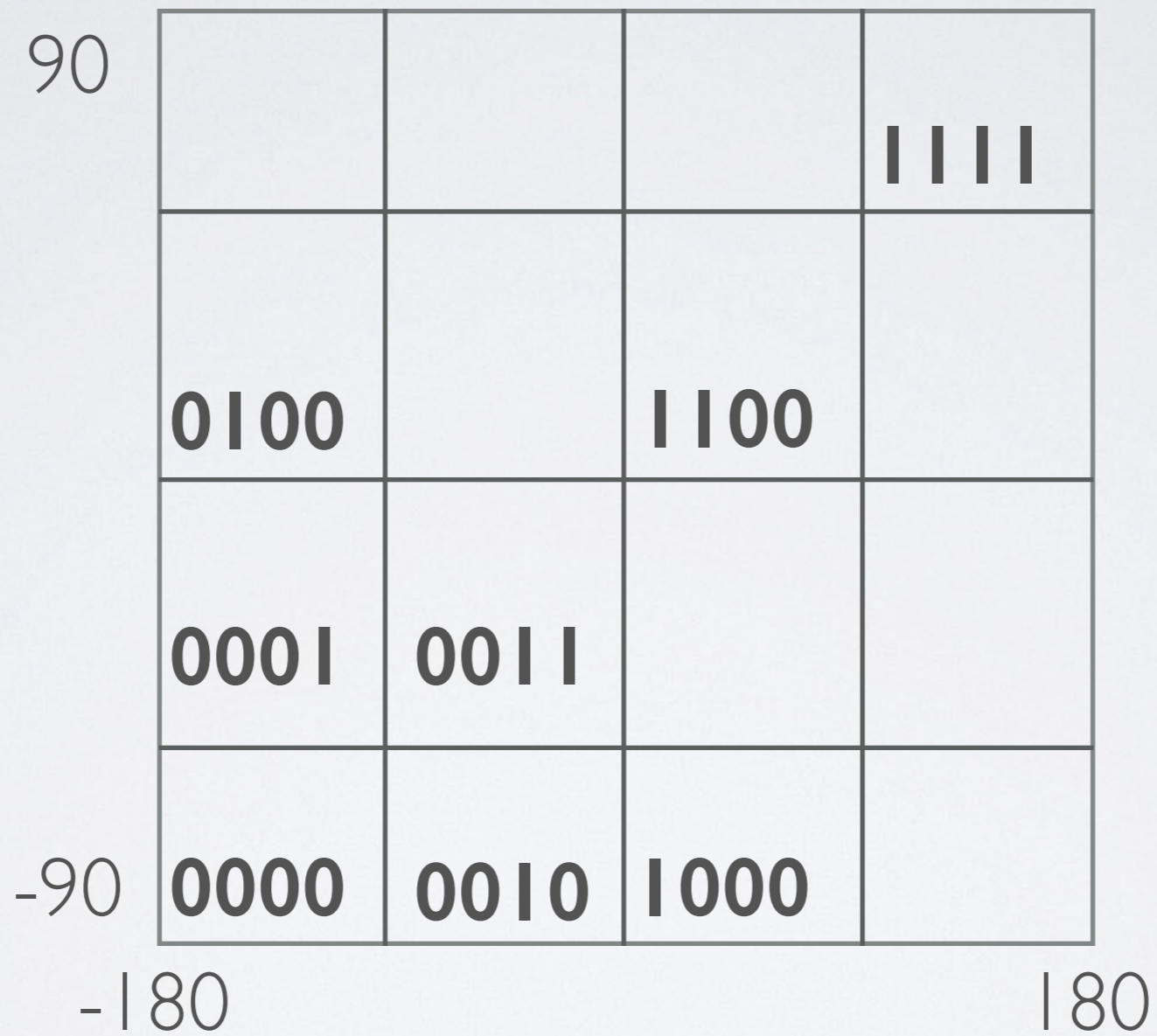
# GEOHASH

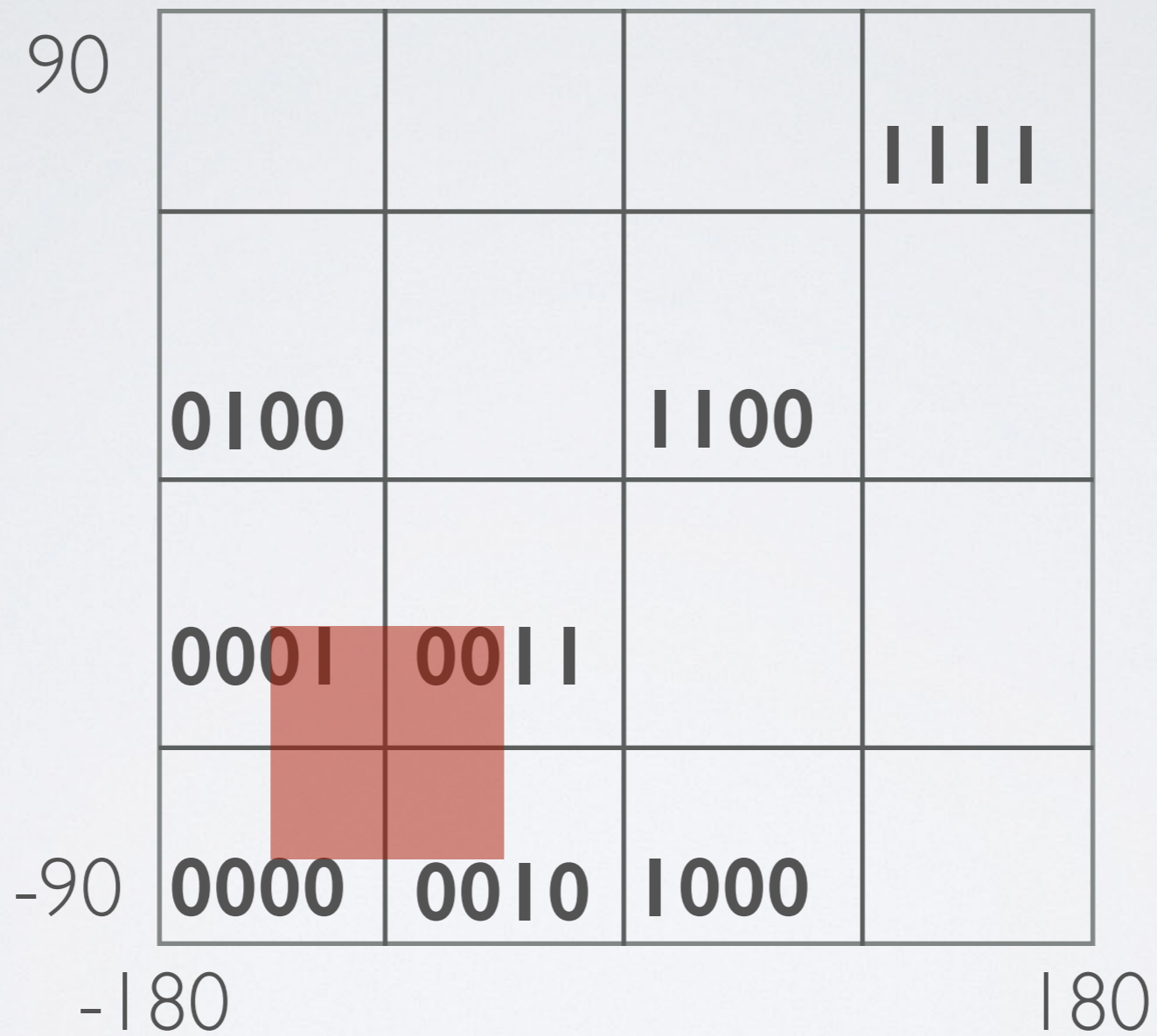
<http://en.wikipedia.org/wiki/Geohash>



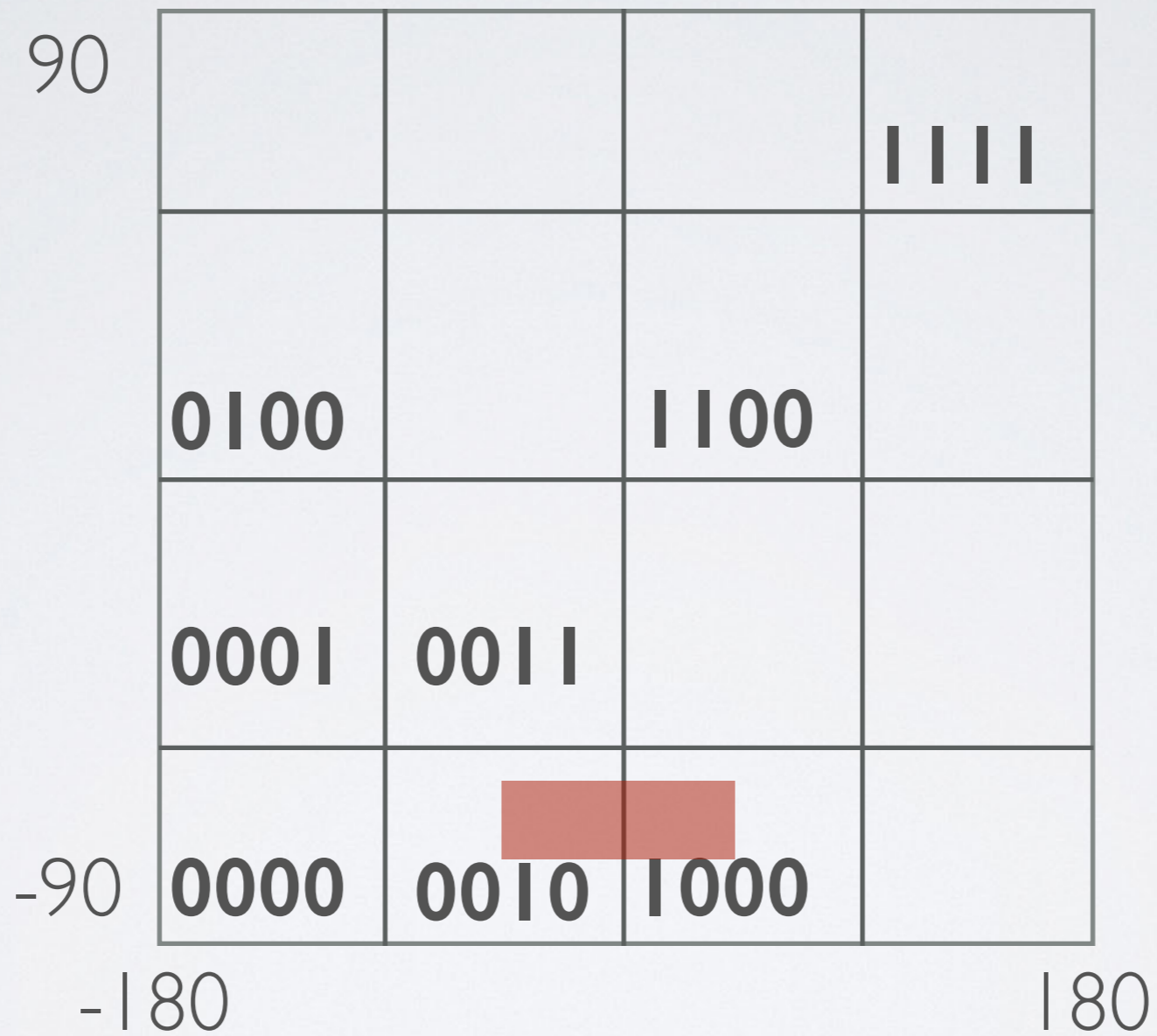
if left of vertical center set left bit to 0 else 1  
if lower of horizontal center set right bit 0 else 1











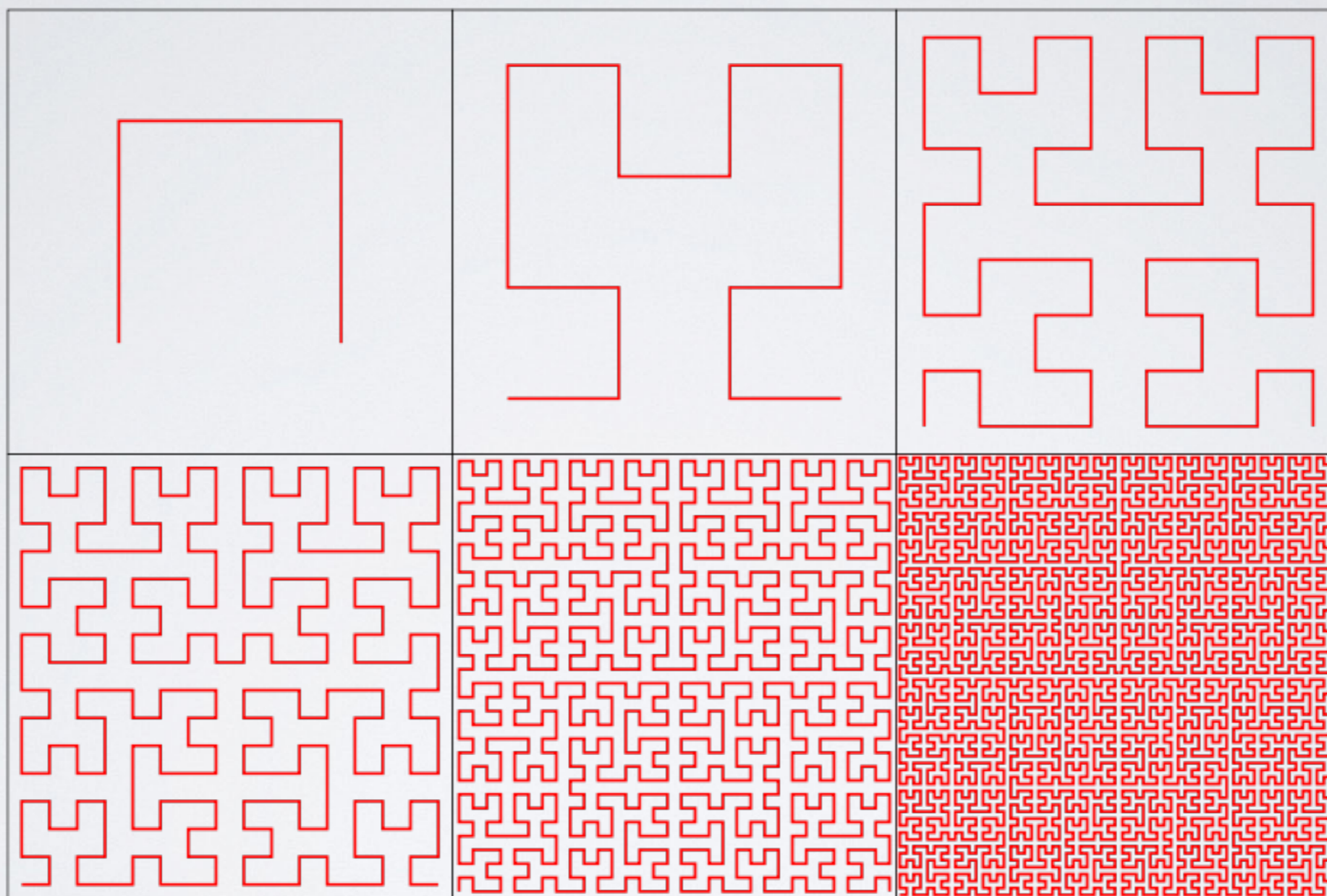
# SPACE FILLING CURVES

~ 1880

N DIM → I DIM

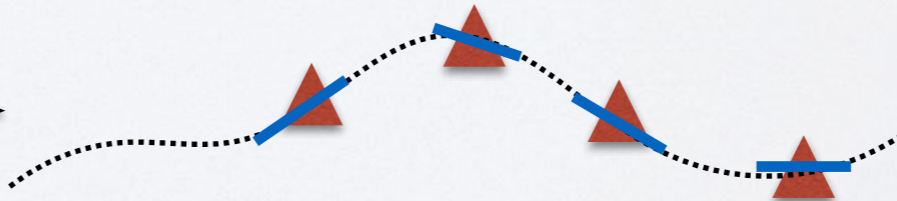
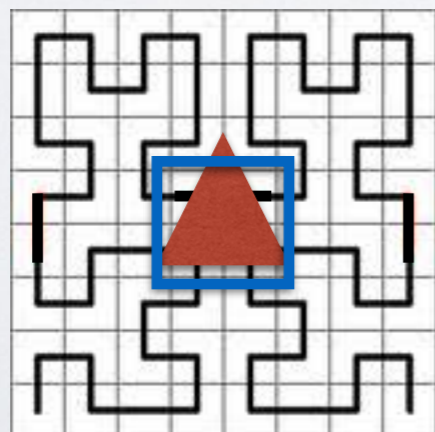
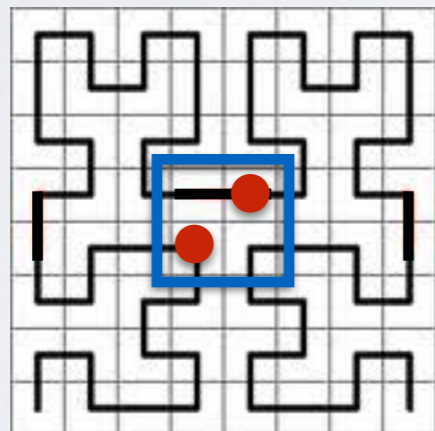
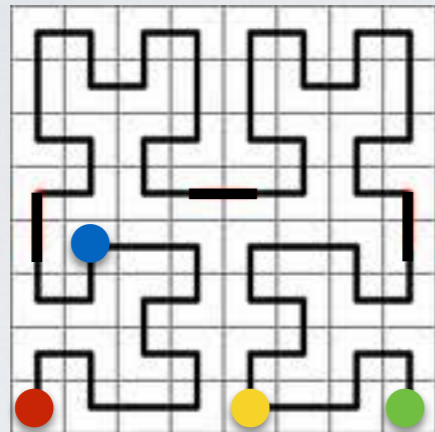


# HILBERT CURVE



[http://en.wikipedia.org/wiki/Space-filling\\_curve](http://en.wikipedia.org/wiki/Space-filling_curve)

# SPACE LINEARIZATION

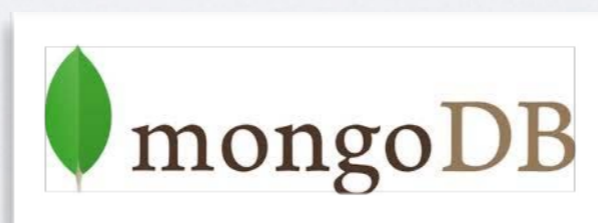




# SPATIAL SUPPORT

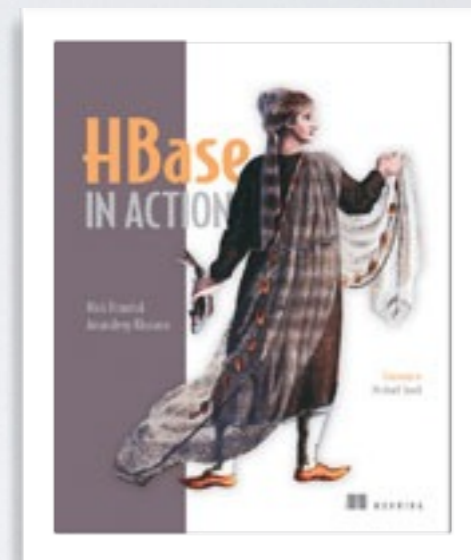
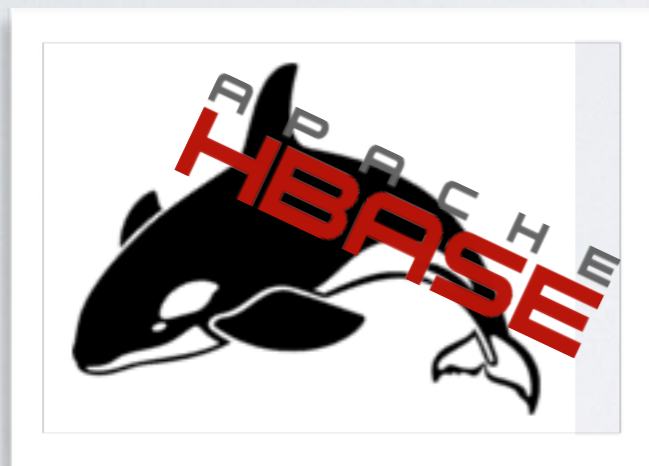


RTree





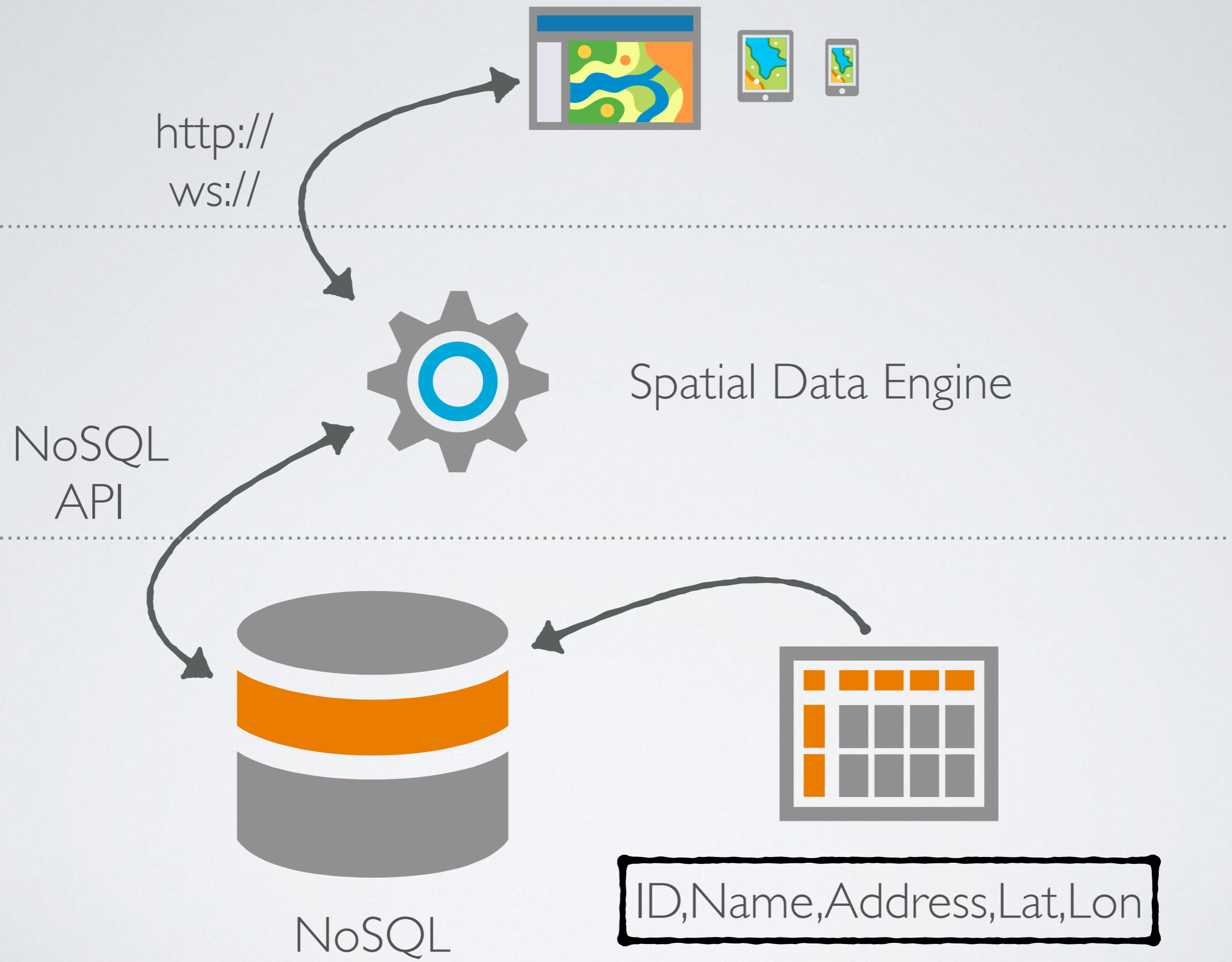
# INDIRECT SUPPORT



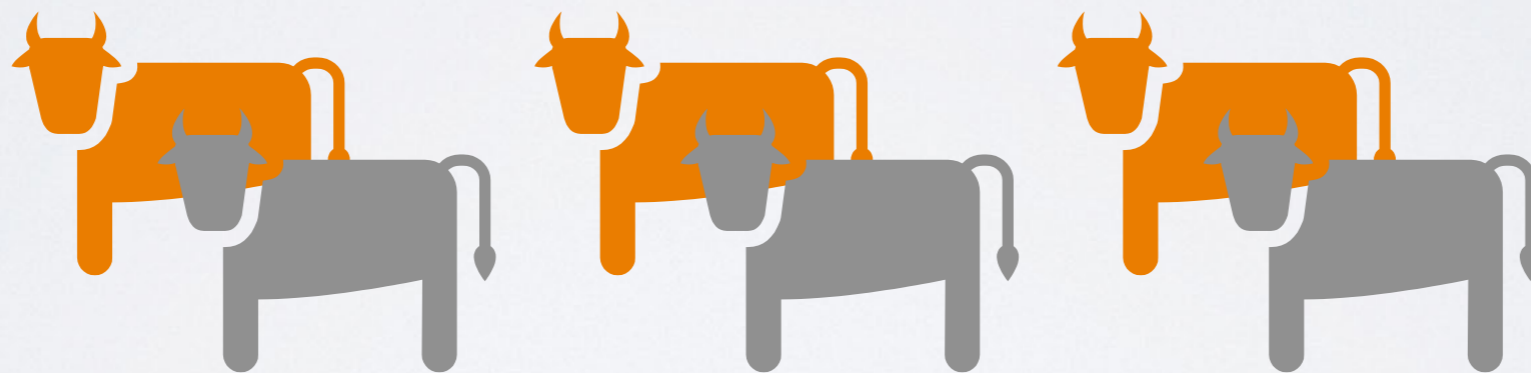
WHAT IS OLD...  
IS NEW AGAIN !

# SPATIAL MIDDLEWARE

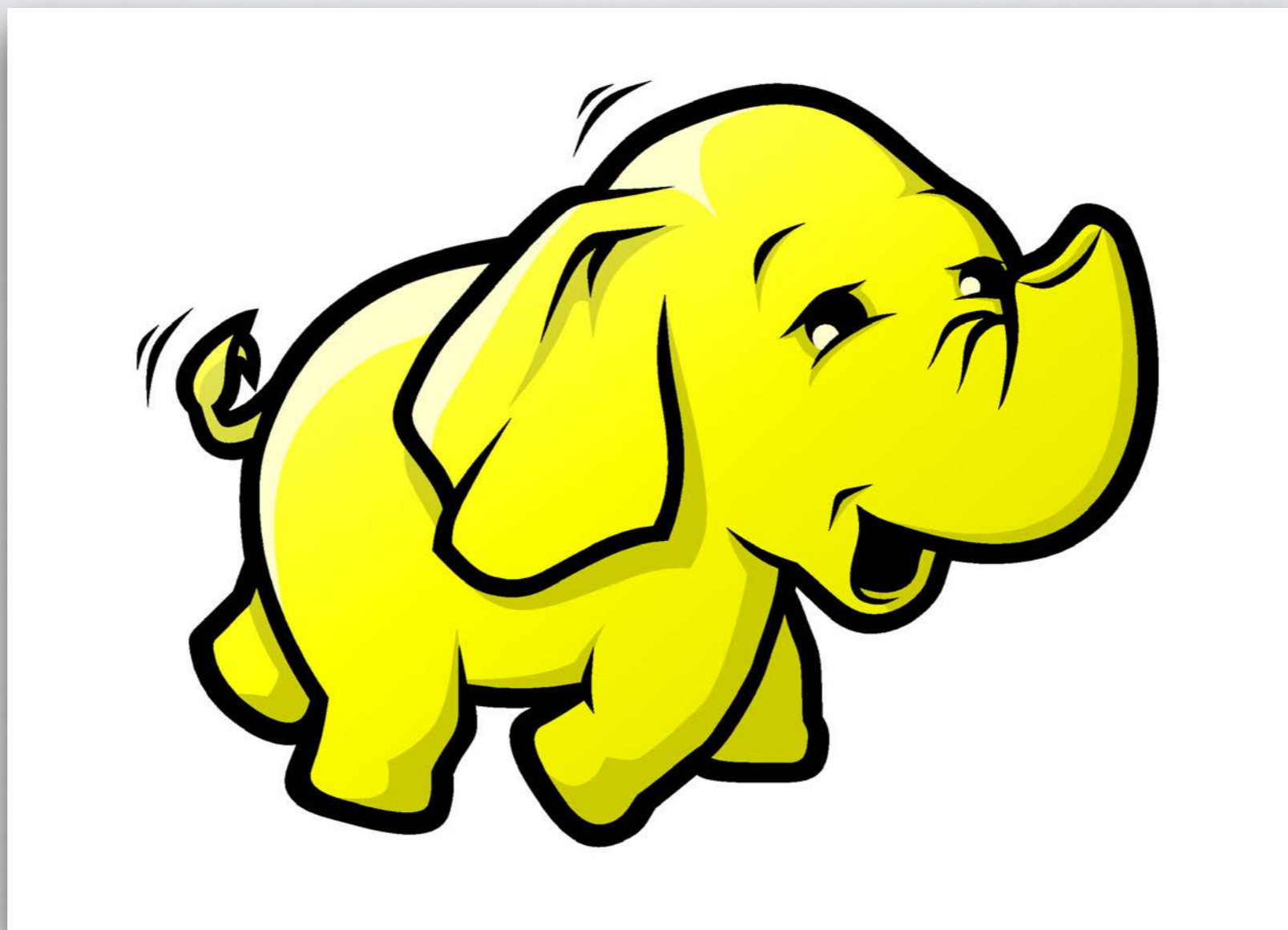




NOT A BIGGER OX...



HADOOP.APACHE.ORG





# WHAT'S IN A NAME ?



<http://blog.pivotal.io/pivotal/products/demystifying-hadoop-in-5-pictures>

# WHAT IS HADOOP ?

- Library / Framework
- Very Very Large Un/Structured Dataset
- Multi Node Distributed Processing
- Resilient To Commodity Hardware Failure



# HADOOP BASIC STACK

MapReduce

Yet Another Resource Negotiator (YARN)

Hadoop Distributed File System (HDFS)





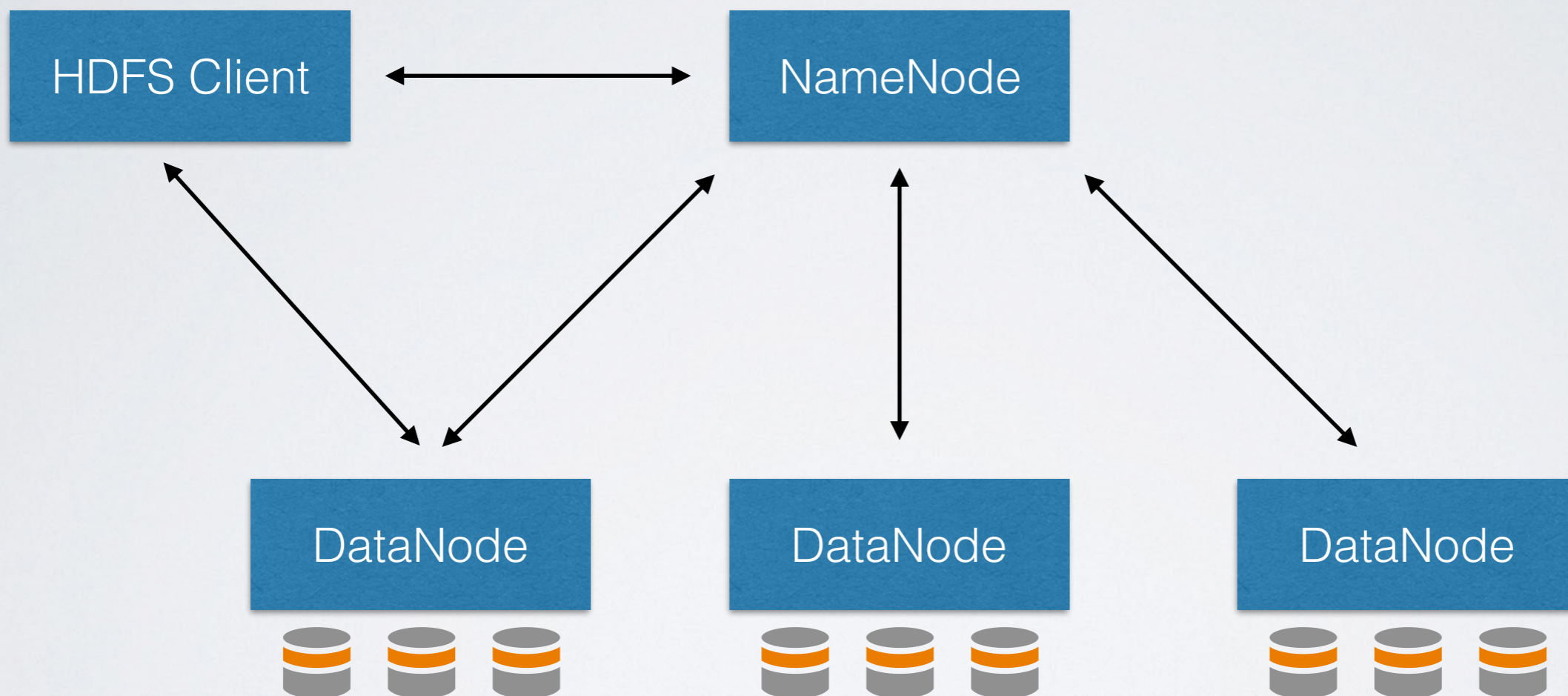
# OTHER HADOOP PROJECTS

- Avro - Serialization / RPC System
- HBase - Distributed Columnar Database
- Hive - Ad Hoc “SQL” Interface
- Pig - Data Flow Parallel Execution (AML)
- ZooKeeper - Coordination Service
- More.....

# HDFS

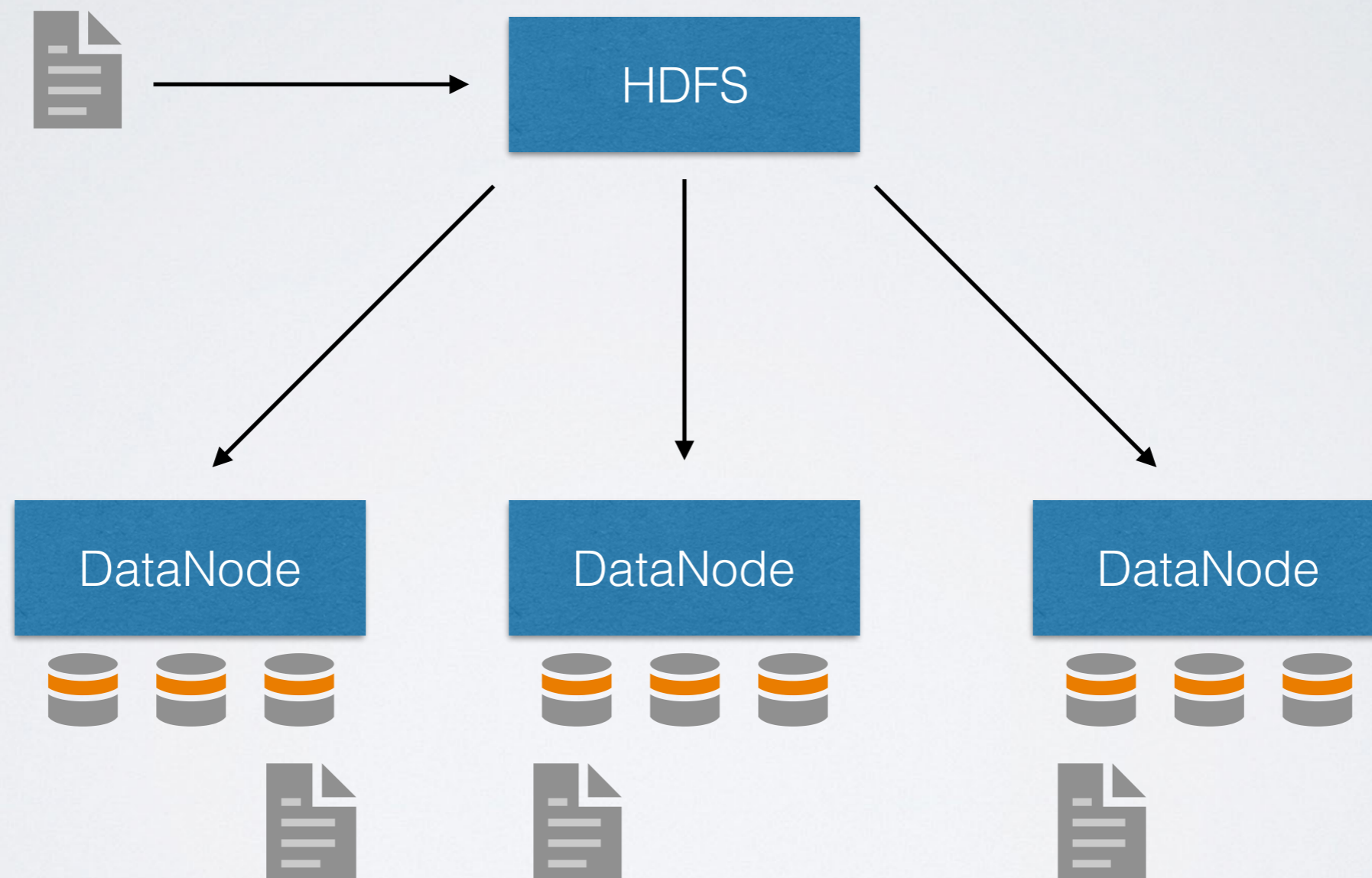
- Distributed File System
- Lots and Lots of Commodity Drives
- Fault Tolerant
- Loves Big Files
- “POSIX” Like Interface

# HDFS





# HDFS Resilience !



# Program



# BigData

# Program



# BigData



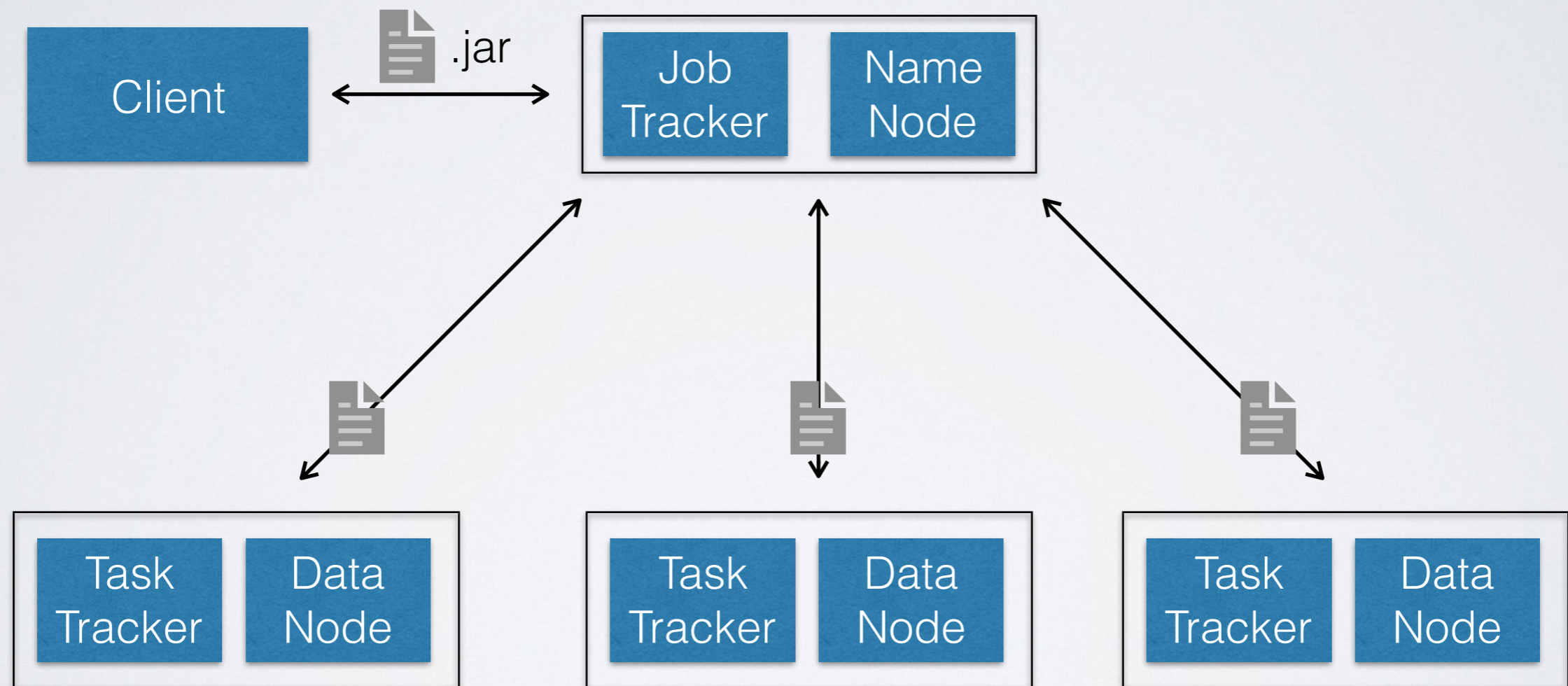
# MAPREDUCE

[http://hadoop.apache.org/docs/r1.2.1/mapred\\_tutorial.html](http://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html)

# WHAT IS MAPREDUCE ?

- Parallel Fault Tolerant Framework
- Splits Large Input
- Invoke User Defined “Map” Function
- Shuffle and Sort
- Invoke User Defined “Reduce” Function

# MAPREDUCE & HDFS





WRITING MR IS HARD...

HOW ABOUT.....  
NO PROGRAMING ???





# APACHE HIVE

“SQL”



MapReduce Job

# HQL

**drop table if exists** *logs*;

**create external table if not exists** *logs*(  
ip **string**,  
method **string**,  
uri **string**,  
status **string**,  
bytes **int**,  
time\_taken **int**,  
referrer **string**,  
user\_agent **string**  
) **partitioned by** (*year int, month int, day int, hour int*)  
**row format delimited**  
**fields terminated by** '\t'  
**lines terminated by** '\n'  
**stored as** *textfile*  
**location** 'hdfs://hadoop:8020/logs/';


# OTHER ADHOC ENGINES

- Cloudera Impala
- Facebook Presto
- SparkSQL
- Bypass MR generation / Direct HDFS Access




WHAT ABOUT SPATIAL ?

GIS Tools for Hadoop by Esri  
GIS Tools for Hadoop by Esri



# GIS Tools for Hadoop

Big Data Spatial Analytics  
for the Hadoop Framework



View project on  
GitHub

Looking at data without location, most of the time seems like looking at just part of a story. Including location and geography in analysis reveals patterns and associations that otherwise are missed. As Big Data emerges as a new frontier for analysis, including location in Big Data is becoming significantly important.

Data that includes location, and that is enhanced with geographic information in a structured form, is often referred to as Spatial Data. Doing Analysis on Spatial data requires an understanding of geometry and operations that can be performed on it. Enabling Hadoop to include spatial data and spatial analysis is the goal of this Esri Open Source effort.

**GIS Tools for Hadoop** is an open source toolkit intended for Big Spatial Data Analytics. The toolkit provides different libraries:

- **Esri Geometry API for Java:** A generic geometry library, can be used to extend Hadoop core with vector geometry types and operations, and enables developers to build MapReduce applications for spatial data.
- **Spatial Framework for Hadoop:** Extends Hive and is based on the

is maintained by **Esri**.

This page was generated by [GitHub Pages](#) using the Architect theme by [Jason Long](#).

Display a menu



# GIS TOOLS FOR HADOOP

- Computational Geometry Library
- Hive Spatial UDF Functions
- GeoProcessing Extensions to ArcMap



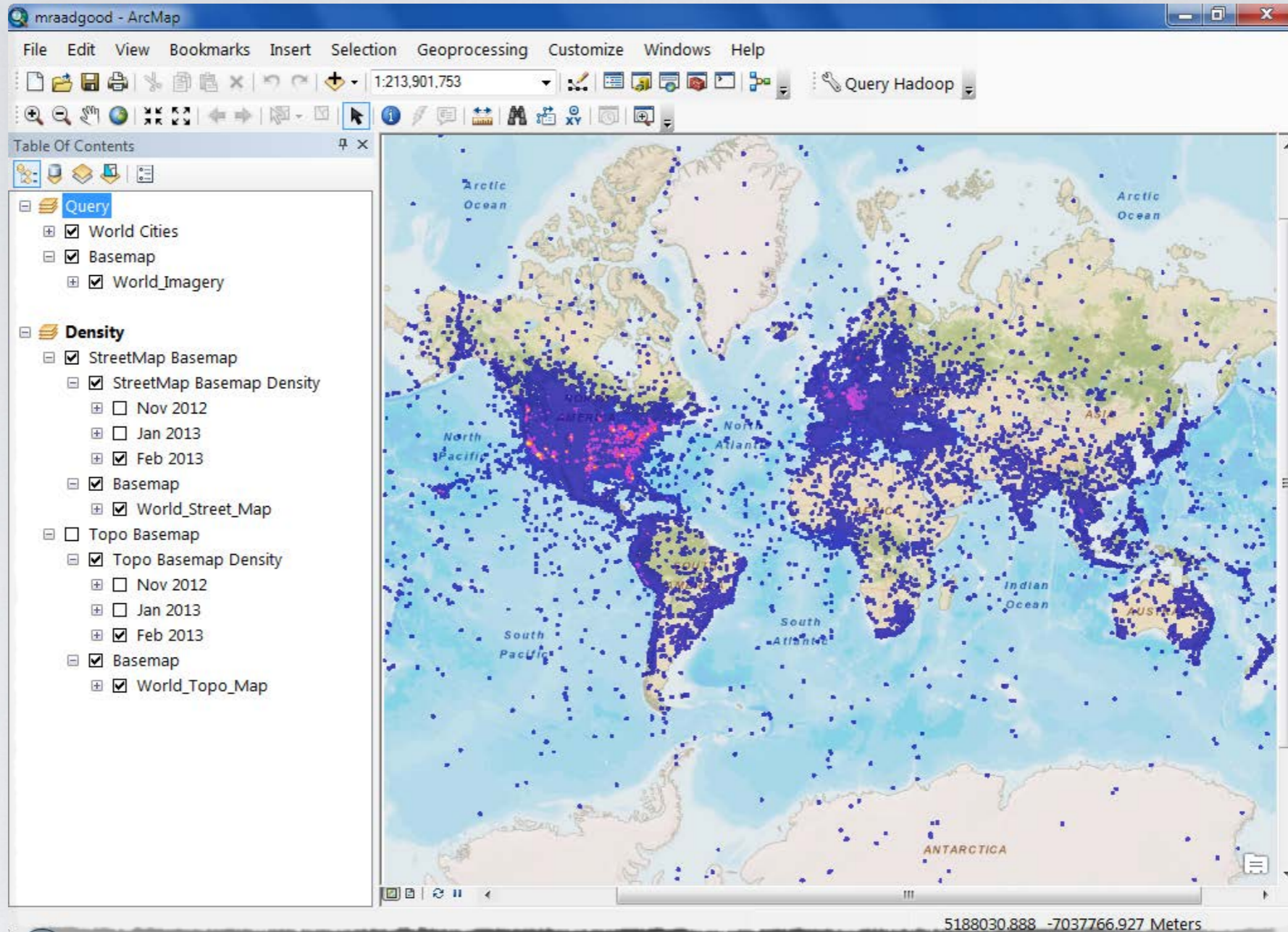
# GEOMETRY LIBRARY

- Points / Lines / Polygons
- I/O (GeoJSON,WTK,WBT,Shape)
- Spatial Relations (inside, touches, intersects,...)
- Spatial Operations (buffer, cut, convex hull,...)
- In-Memory Spatial Index

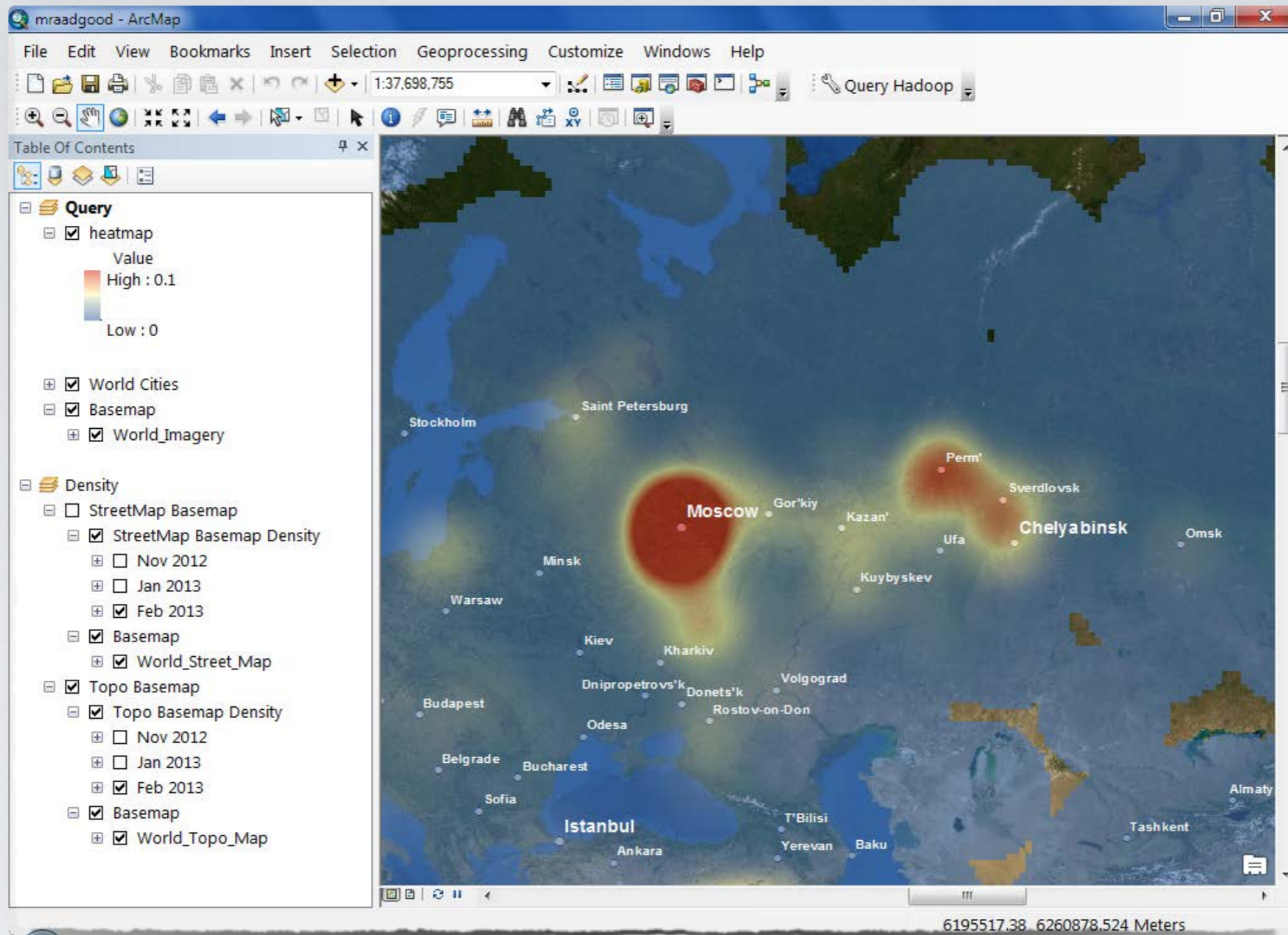
# API USAGE IN BIGDATA

- Map-only jobs - GeoEnrichment
  - Given set of locations
  - Given demographic area
  - Augment location with demographic attributes

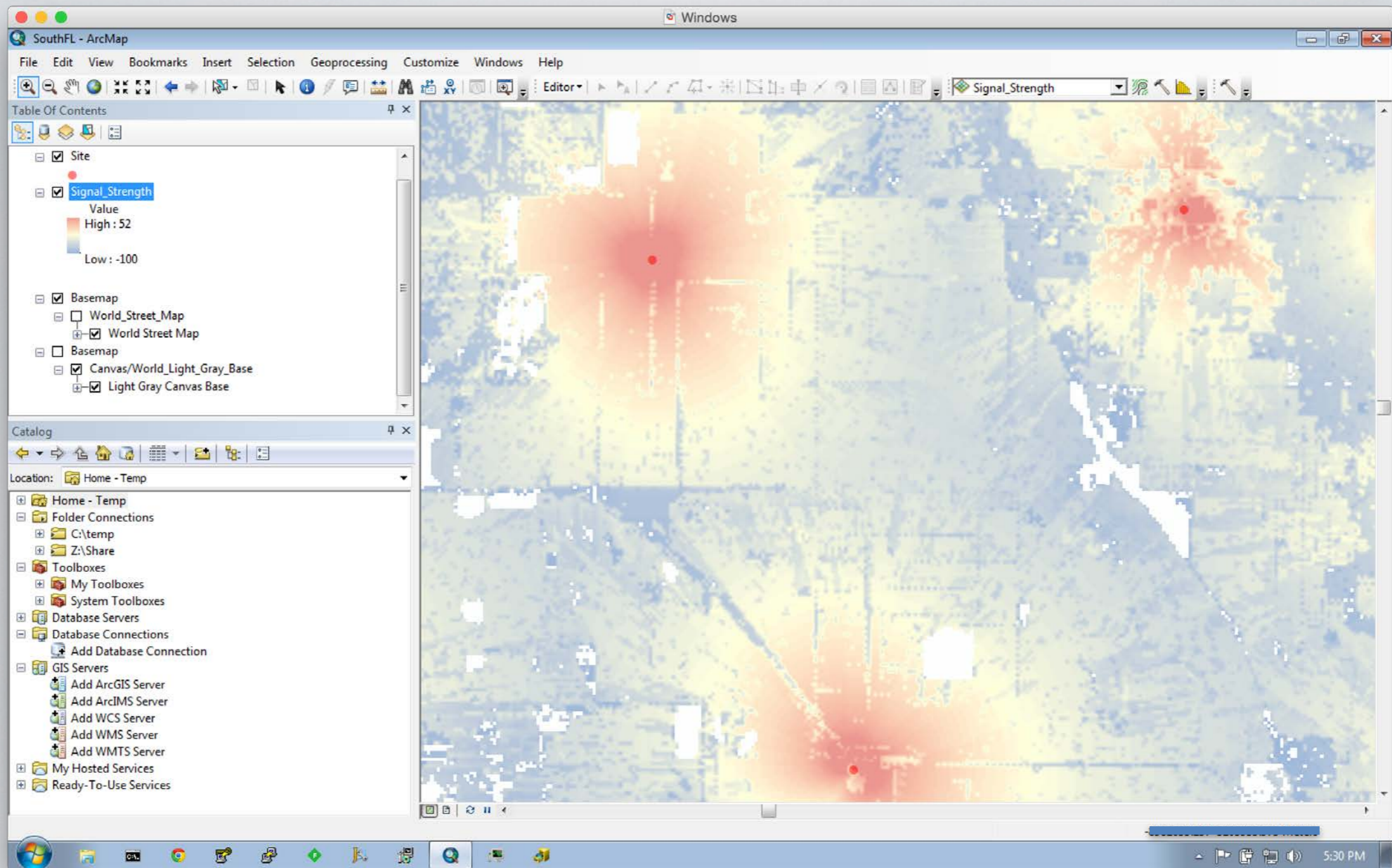












# TELCO



- CDR - Call Data Record
  - DateTime, UUID, LatLon, Duration, Status, etc...
- Drop Call Emerging HotSpot
- Street Traffic Condition
- Massive Spatial Join million x million polygon
- Overlay with Demographic polygons
- Overlay with Current Weather
- Overlay with Social Media



# TELEMATICS

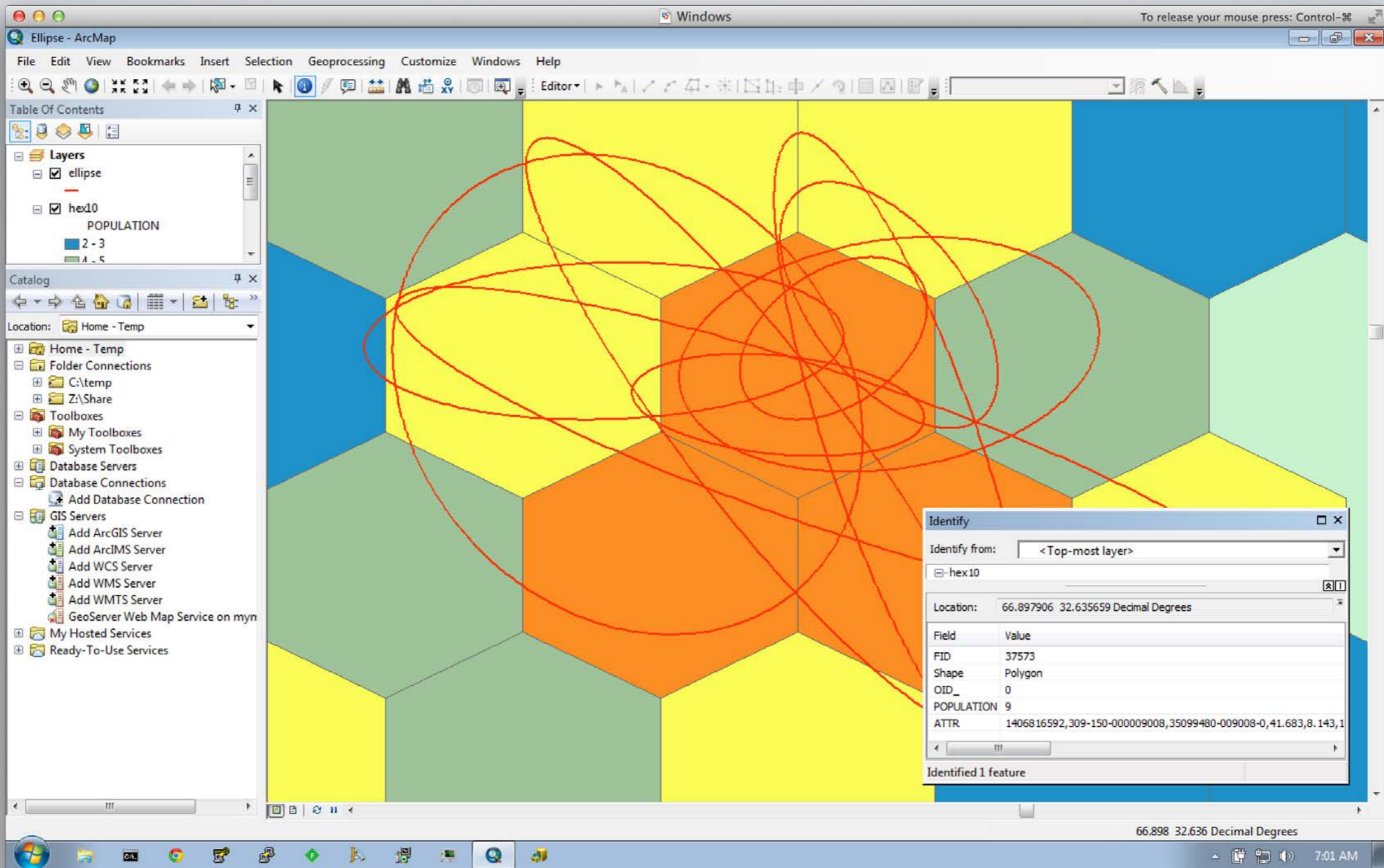
- Feedback to Engineering
- Car to Car Communication
- Street Condition Detection
- Best Route Prediction (EV)
- Overlay with weather



**Insurance as you drive !**

# Observations









# A Hadoop-enabled Ship Tracking Application for the Port of Rotterdam

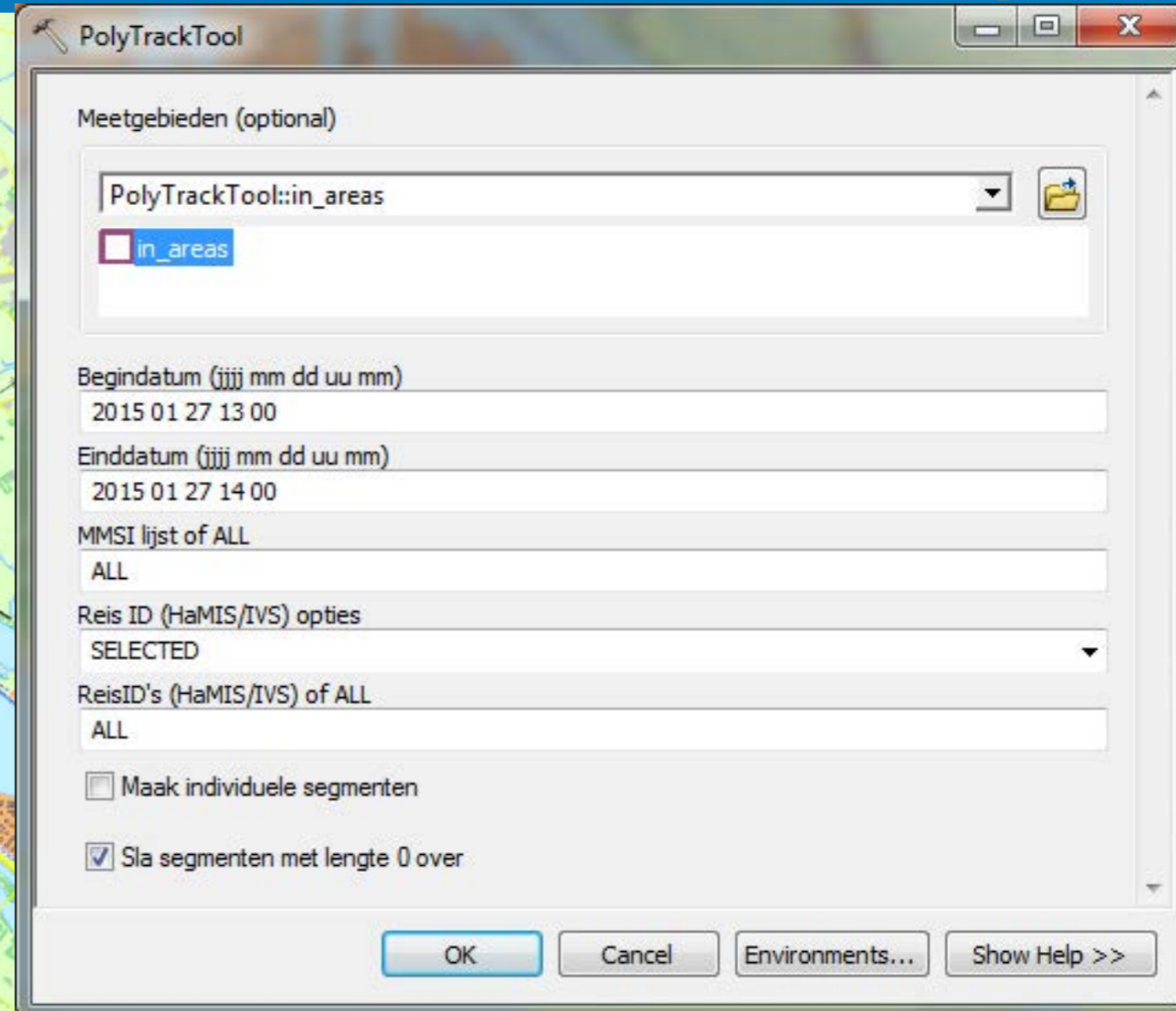
Hadoop Summit, Brussels, 15 April 2015

**Frank Cremer (Geomatik)**

**Mansour Raad (ESRI)**



# Access information in three clicks



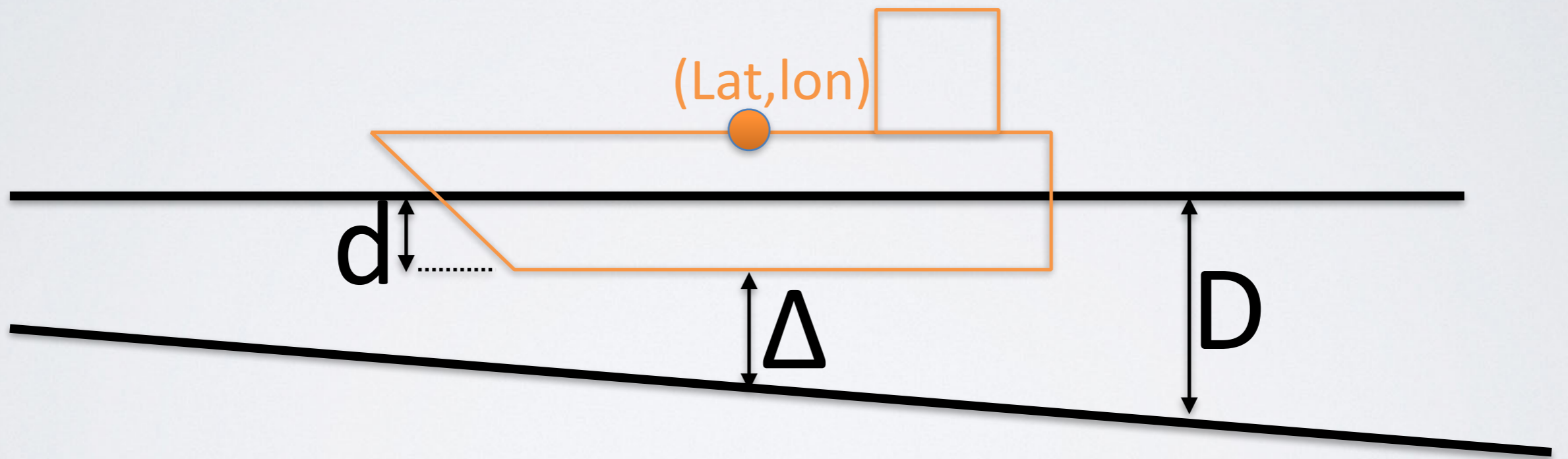


# Usage of ship position data

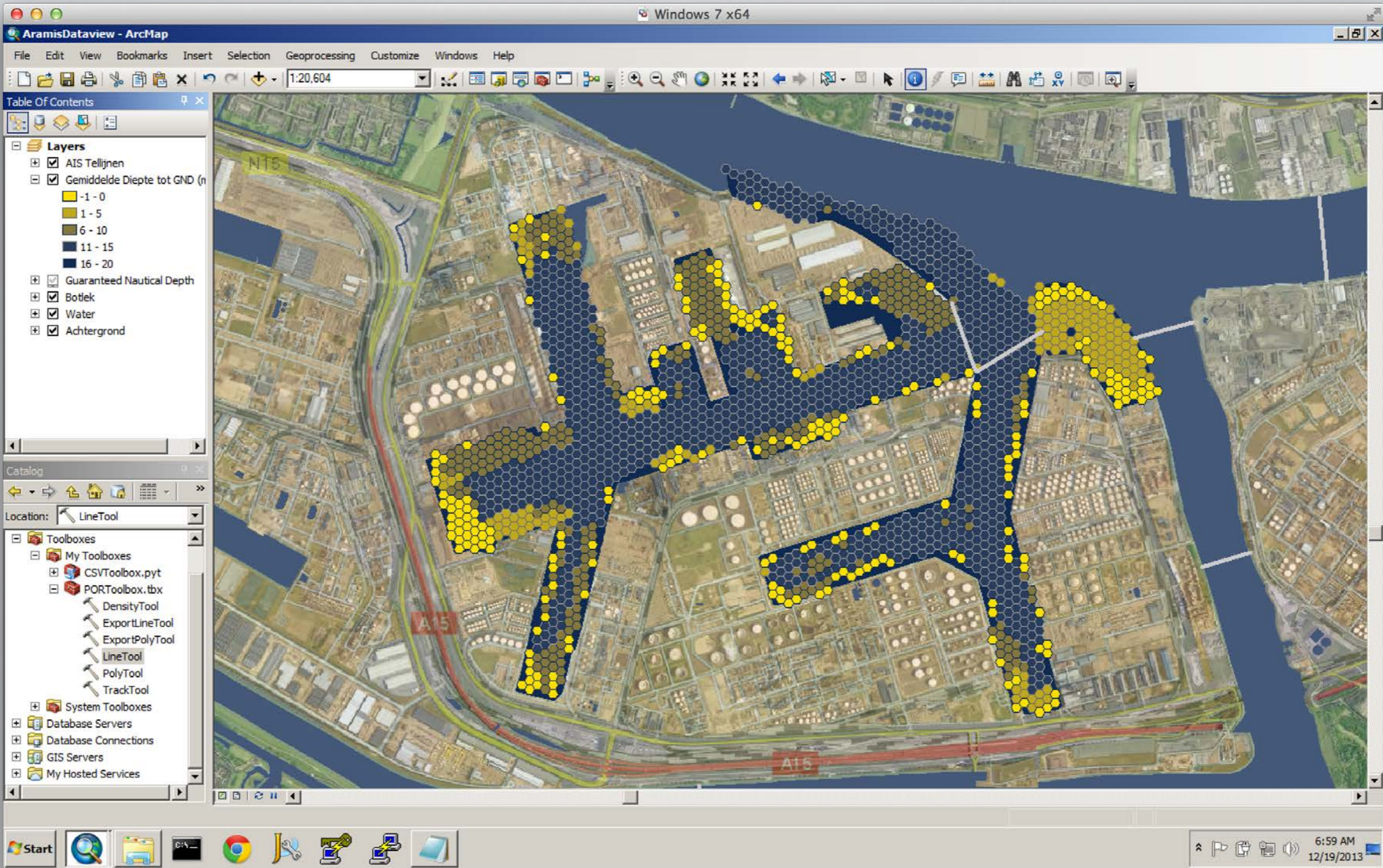
- **Harbour master**
  - Incident analysis
  - Safety checks
- **Capacity management**
  - Identifying bottlenecks
  - Planning decision support
- **Environmental management**
  - Pollution (NOx) calculations
  - Speed measures to reduce pollutions



# Where is $\Delta \approx 0$ ?



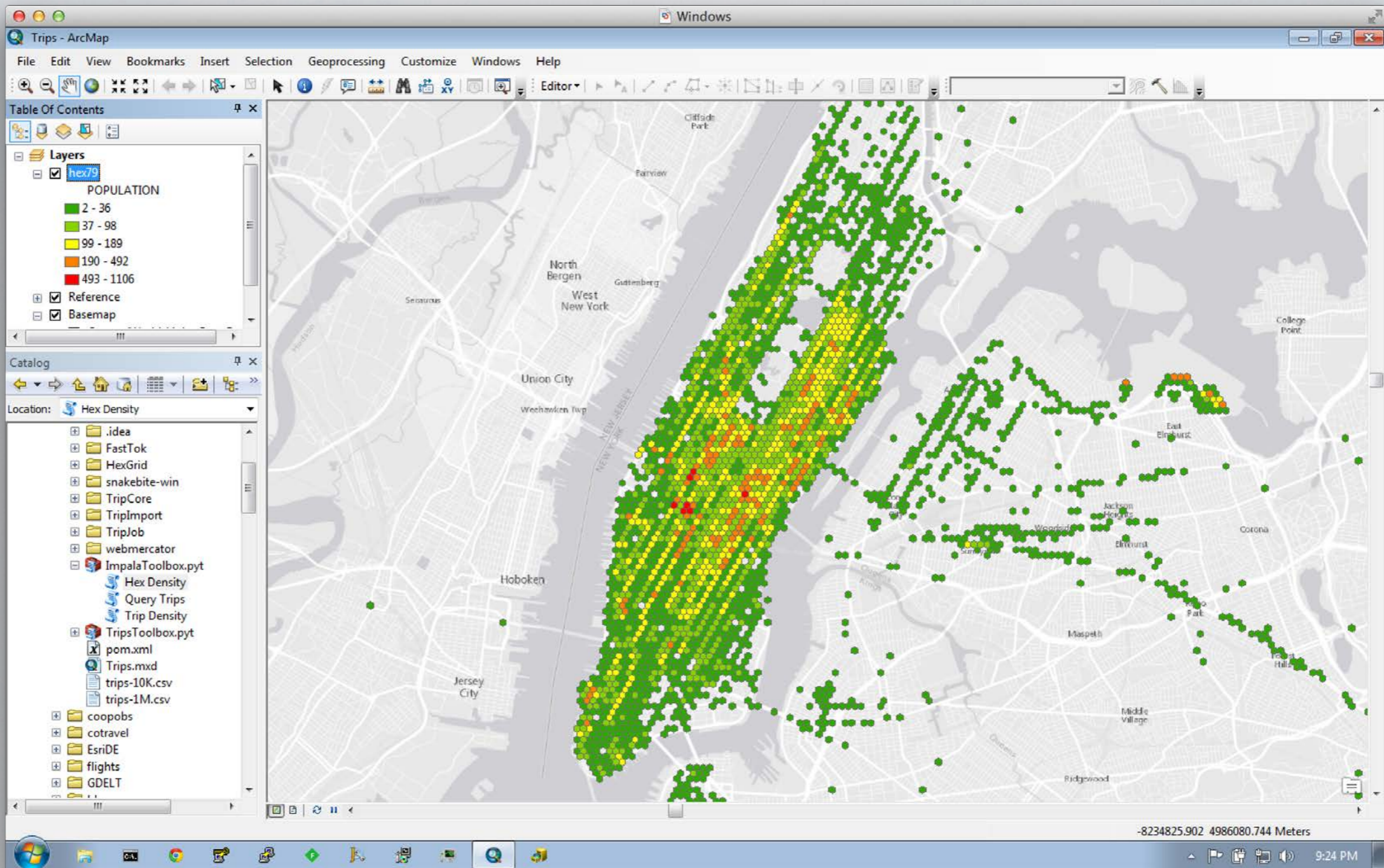






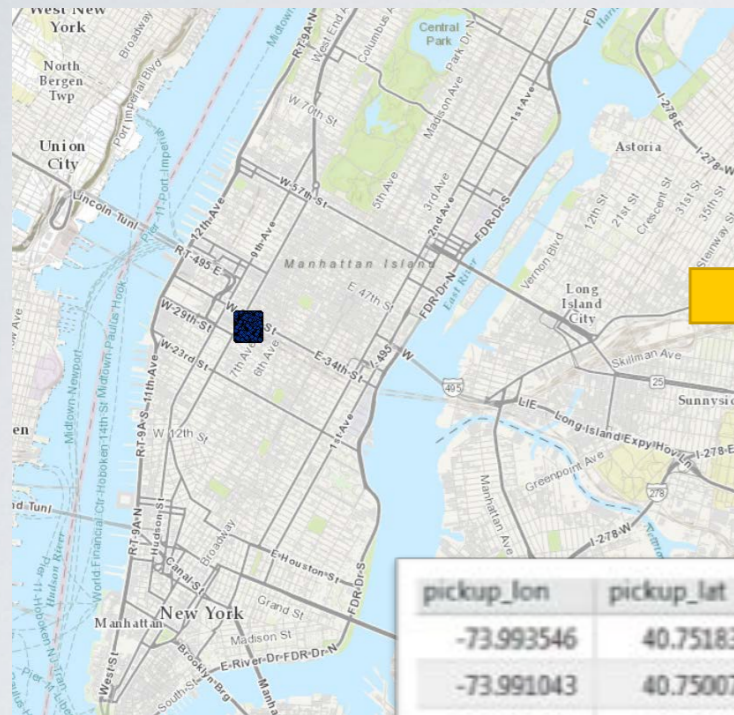






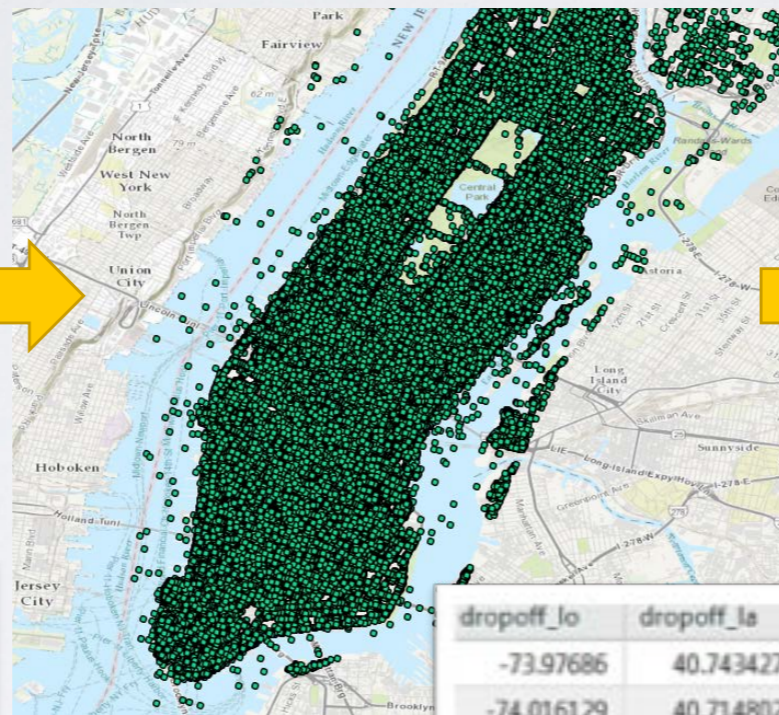


Selected pickups



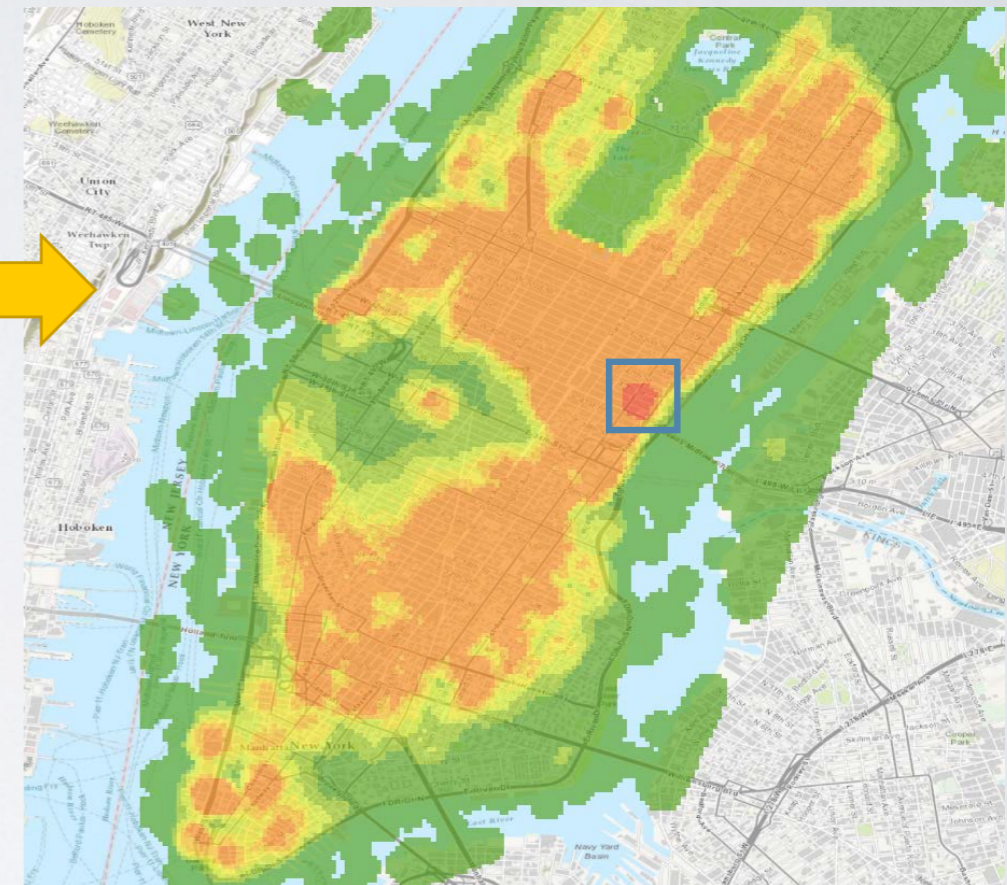
pickup_lon	pickup_lat
-73.993546	40.751839
-73.991043	40.750076
-73.994041	40.751087
-73.991463	40.750256

Corresponding drop-offs



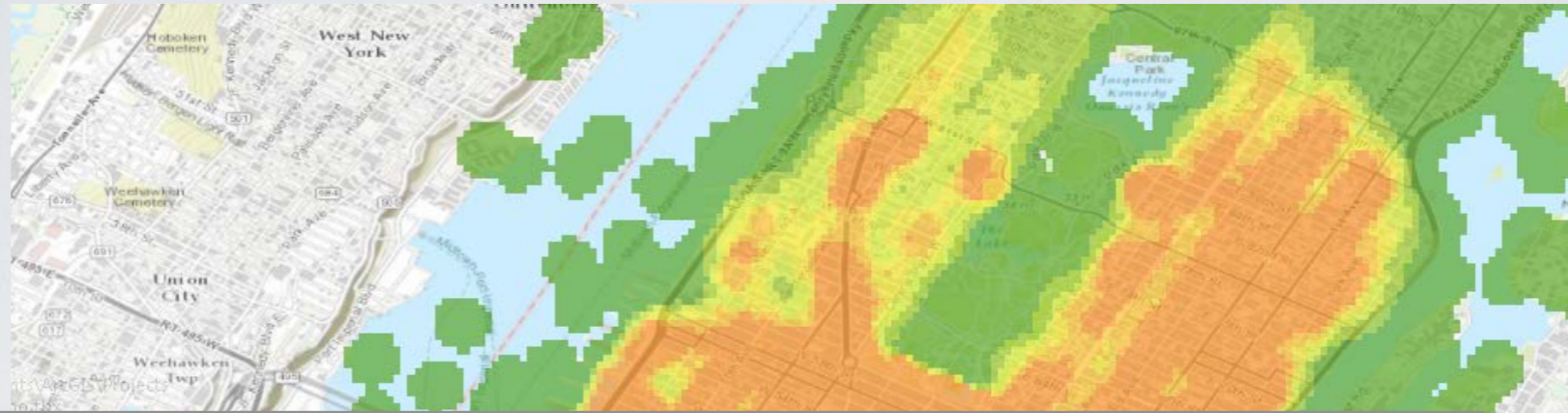
dropoff_lo	dropoff_la
-73.97686	40.743427
-74.016129	40.714802
-73.982033	40.739262
-73.980934	40.730453

Density of passenger drop-offs



Turtle Bay – UN





### Shuttle Locations

- ▲ Pickup
- Drop-off

HOME ▾ My Map NEW MAP

[Details](#)
[Add](#)
[Basemap](#)
[Save](#)
[Share](#)
[Print](#)
[Directions](#)
[Measure](#)
[Bookmarks](#)

**Directions**

A NJ TRANSIT-Penn Station-New Yo  
B United Nations Plz, New York, Nev

[ADD DESTINATION](#)  
[BY CAR](#) [WALKING](#)  
[SHOW MORE OPTIONS](#)  
[GET DIRECTIONS](#) [ADD AS LAYER](#)  
[CLEAR](#)

**1.86 miles · 11 minutes**

[ZOOM TO FULL ROUTE](#)

- A 1. Start at NJ TRANSIT-Penn Station-New York
- 2. Go northwest on **W 31st St** (Joe Louis Plz)  
0.12 mi 1 minute
- ▲ 3. Turn right on **8th Ave** (Joe Louis Plz)

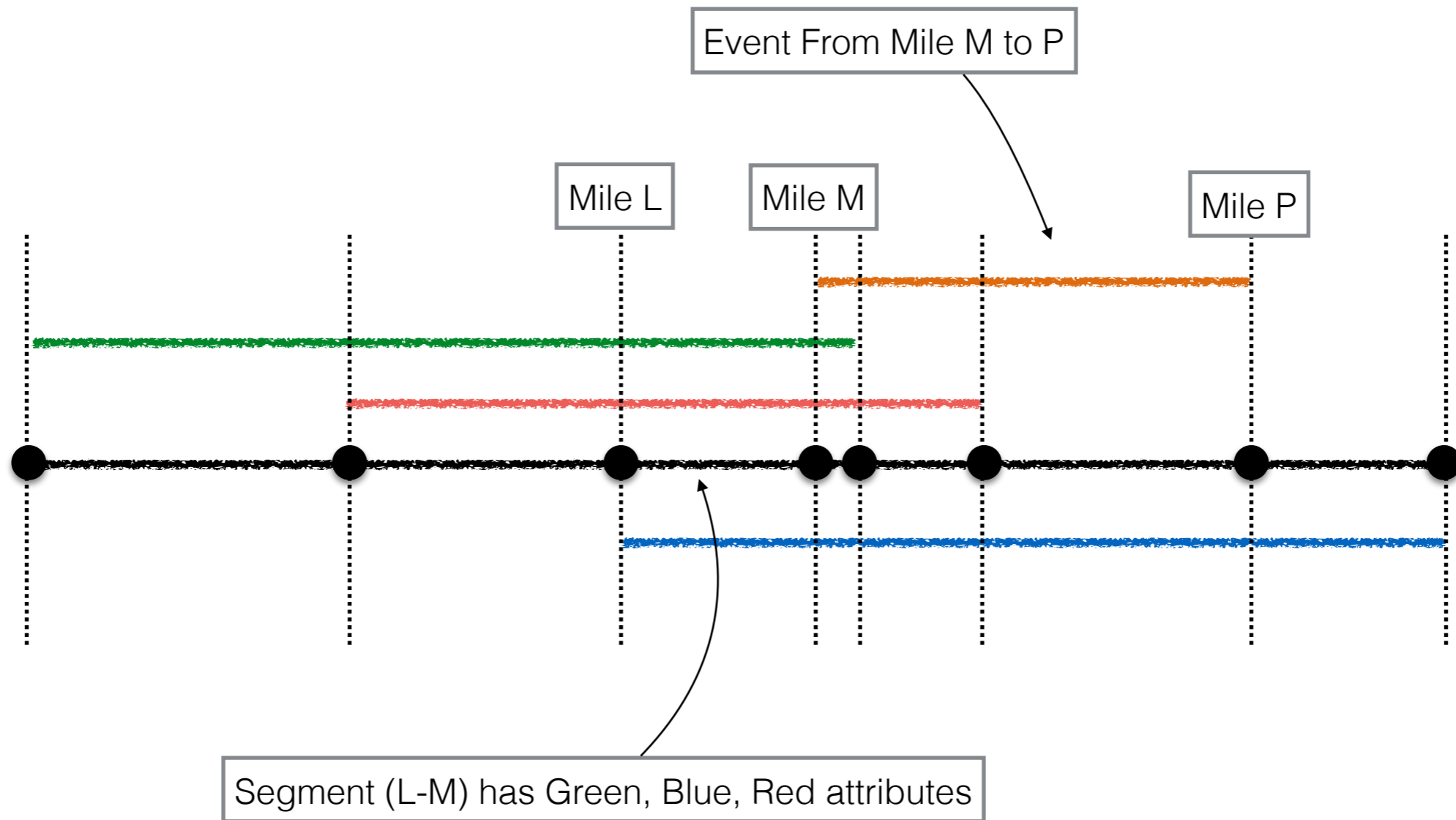




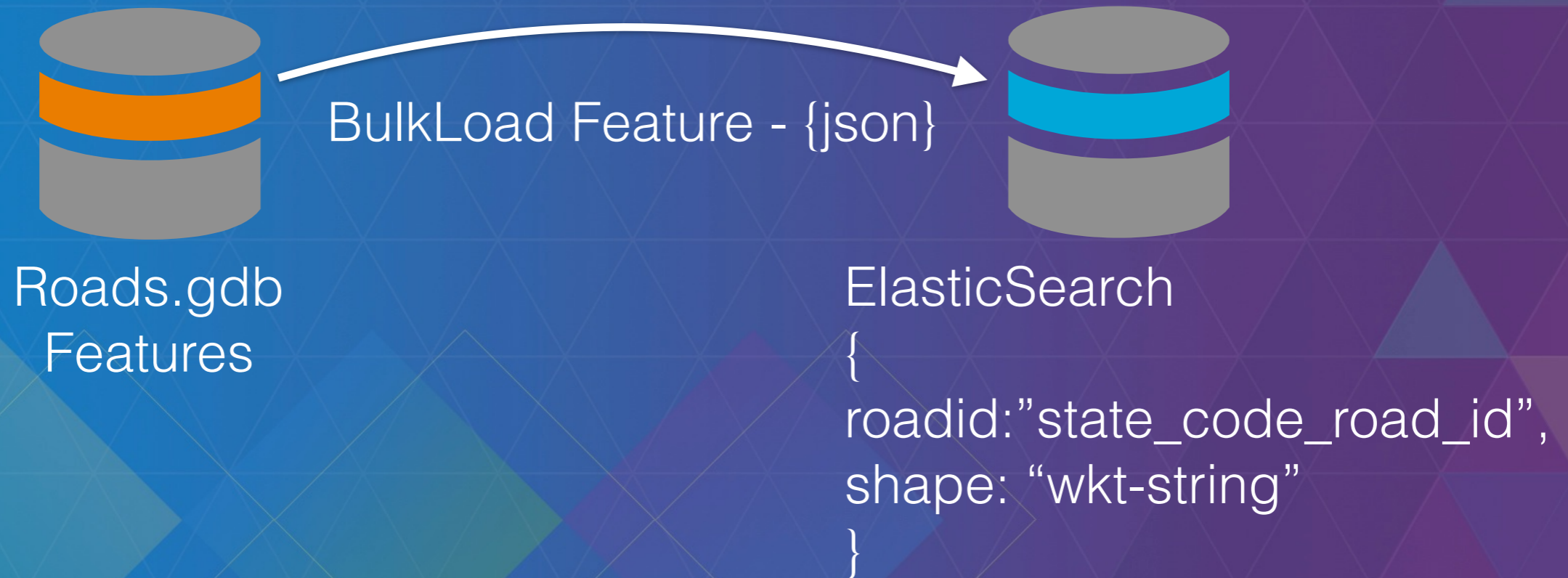


***Mission:* Fast Dynamic  
Segmentation for Linear  
Referencing**





# Step 1 - Clean and Bulk Load Roads





# Step 2 - New Dyn-Seg GeoProcess



## 10M Events

year  
statecode  
routeid  
mile1  
mile2  
key  
val



Parallel  
Distributed  
Share-Nothing  
"GeoProcessing"

{json}

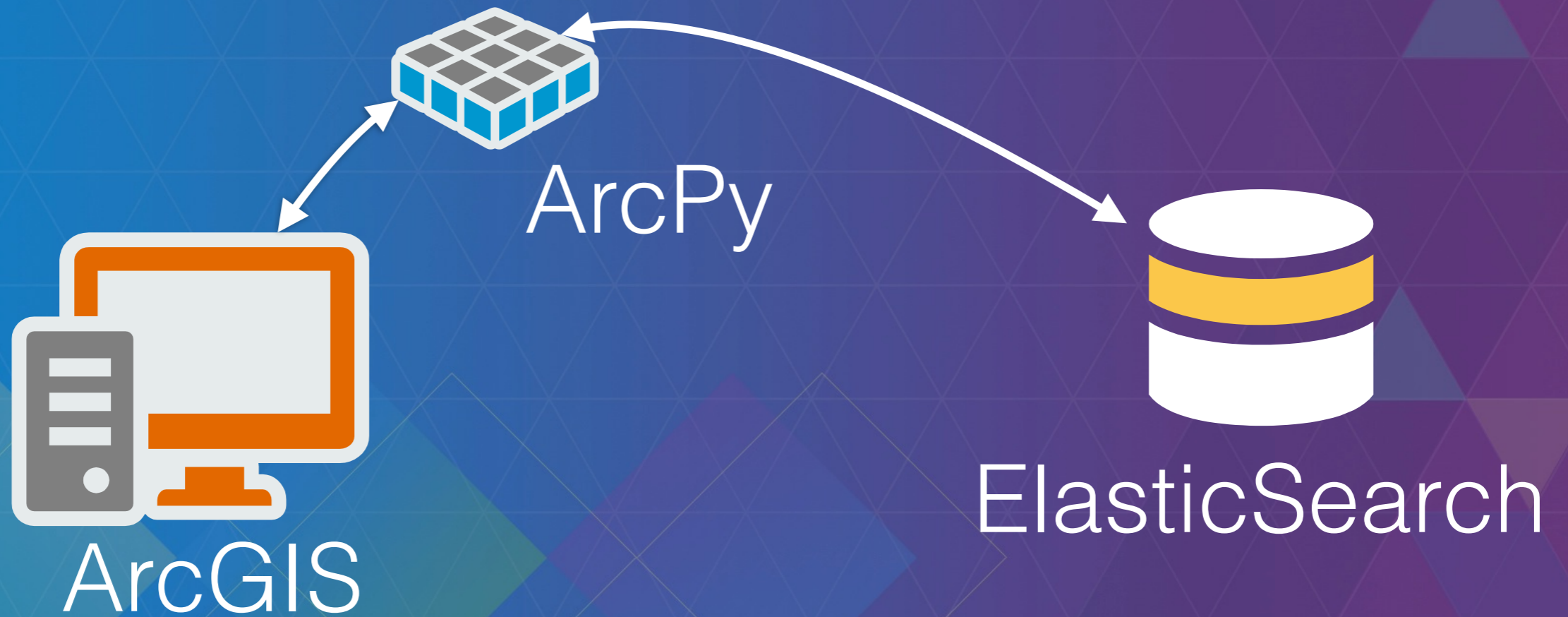


## ElasticSearch

*Index Every Document/Field*

```
{  
  roadid="xxx"  
  wkt="multilinestring(((x y,...)  
  IRI=1  
  F_SYSTEM=2  
  ....  
}
```

## Step 3 - Query and Display



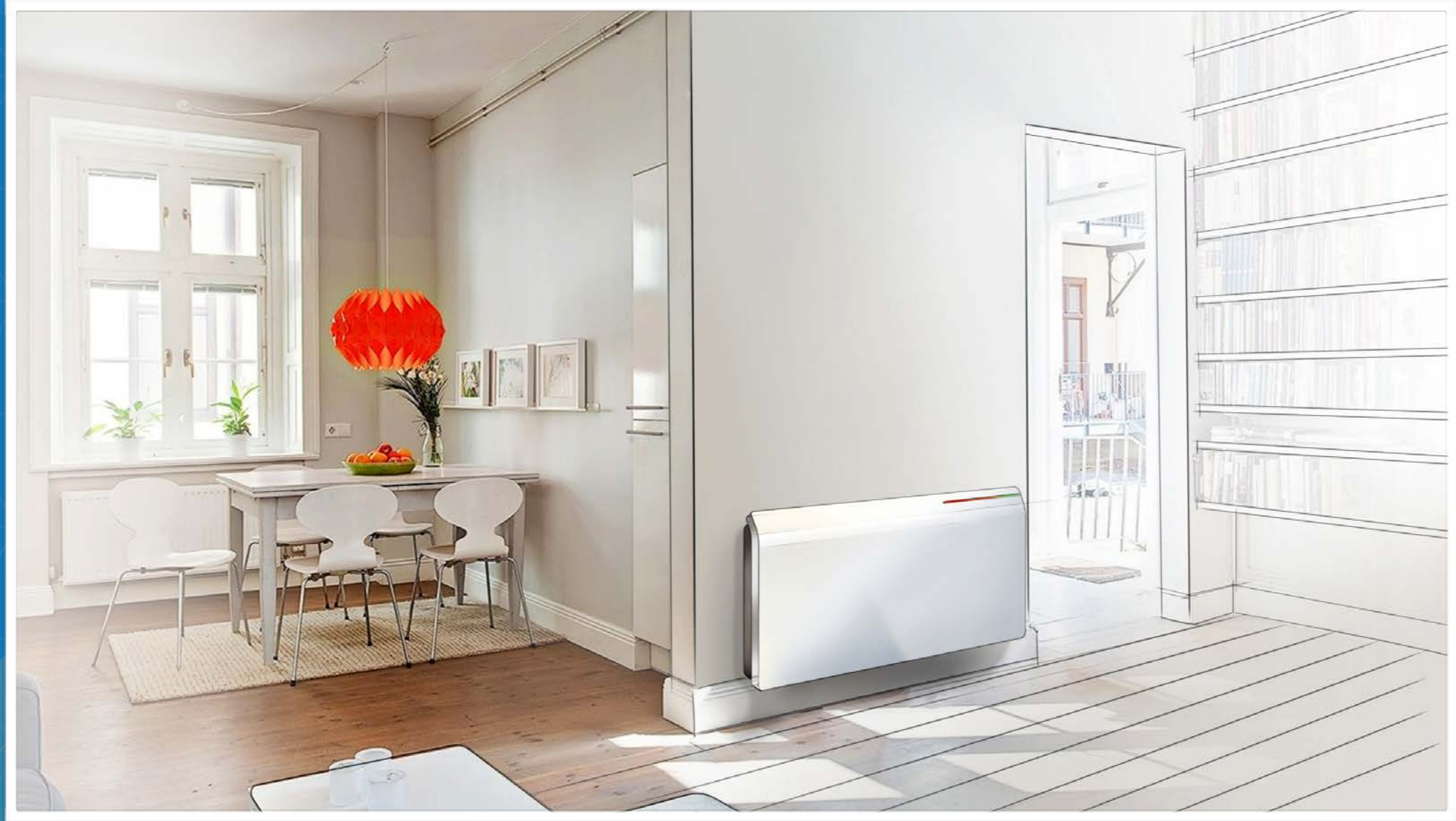


# DataCenters

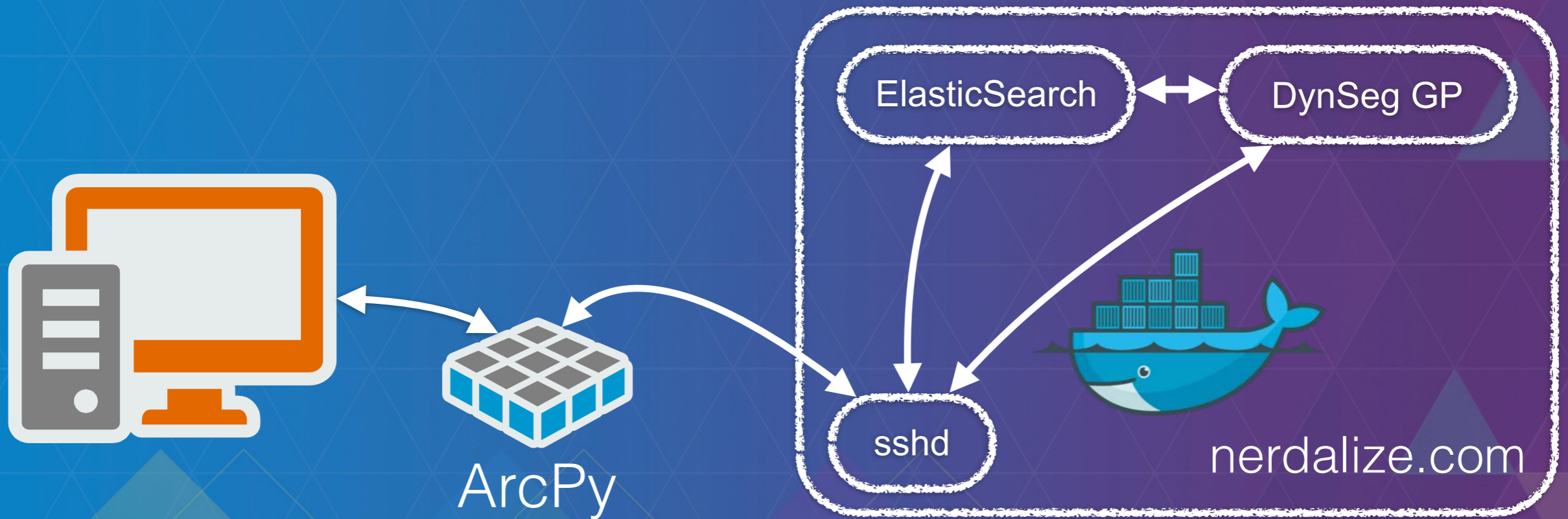




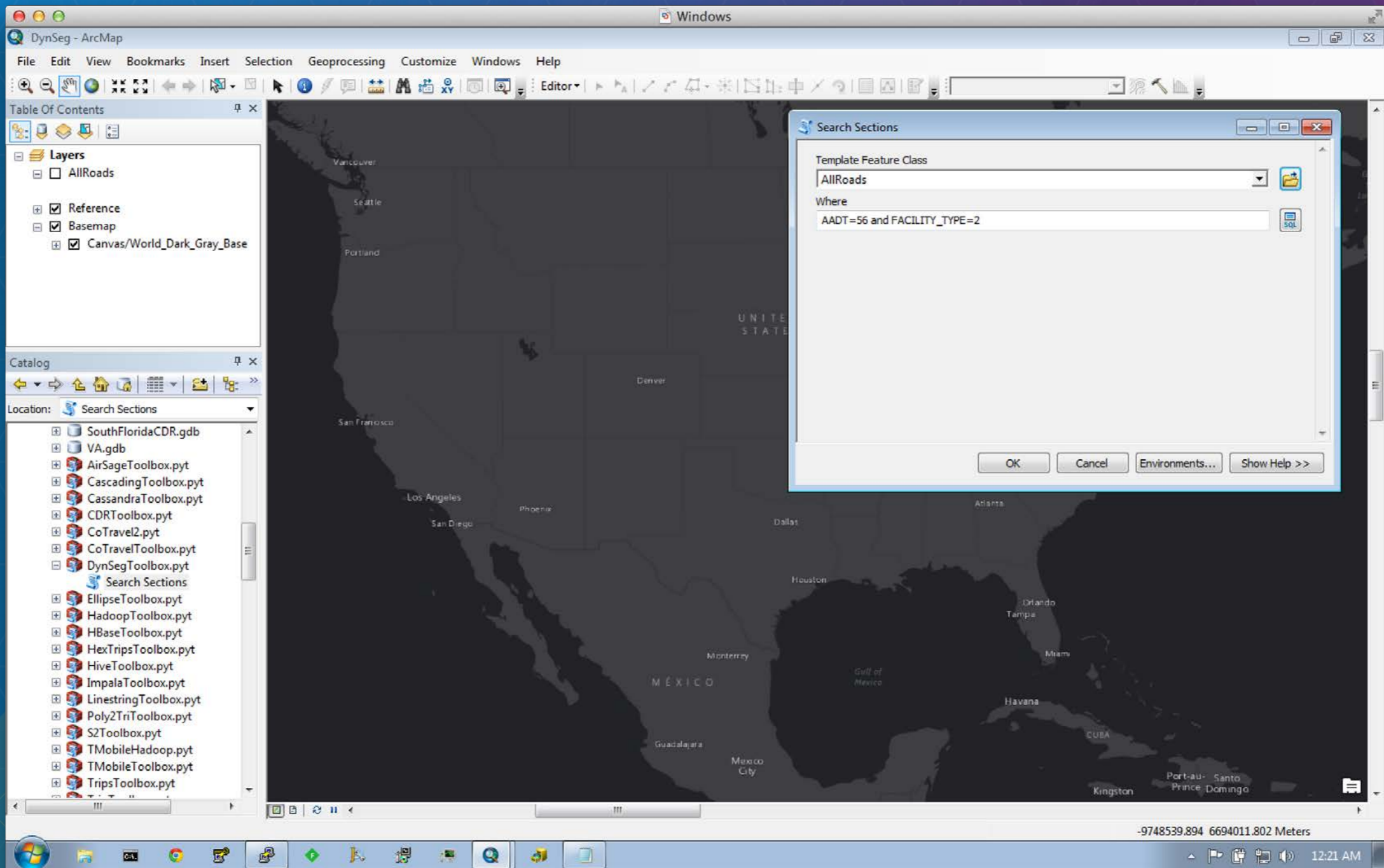




# ArcGIS → Nerdalize







DynSeg - ArcMap

File Edit View Bookmarks Insert Selection Geoprocessing Customize Windows Help

Table Of Contents

- Layers
  - Sections
  - AllRoads
  - Reference
  - Basemap
    - Canvas/World\_Dark\_Gray\_Base

Catalog

Location: Search Sections

- CoTravel2.pyt
- CoTravelToolbox.pyt
- DynSegToolbox.pyt
- Search Sections
- EllipseToolbox.pyt
- HadoopToolbox.pyt
- HBaseToolbox.pyt
- HexTripsToolbox.pyt

Table

Sections

OID	Shape	SHAPE_KEY	YEAR	URBAN_CODE	FACILITY_TYPE	F_SYSTEM	IRI	OWNERSHIP	THROUGH_LANES	AADT
1	Polyline	30_C222319N	2010	99999	2	7	99	2	2	56
2	Polyline	30_C069277N	2010	99999	2	7	99	4	1	56
3	Polyline	30_C015442N	2010	99999	2	7	99	64	2	56
4	Polyline	30_C083210N	2010	99999	2	7	99	60	1	56
5	Polyline	30_C227879N	2010	99999	2	7	99	2	1	56
6	Polyline	30_C247821N	2010	99999	2	7	99	2	1	56
7	Polyline	E1 D V A 0 8 0 5 C 0 0 8 7 1 M R	2010	75421	2	7	99	1	2	56

(0 out of 199 Selected)

Sections

-12724773.814 6221929.574 Meters

12:24 AM



YOU CAN DO IT TOO !





www.cloudera.com/downloads/quickstart\_vms/5-5.5

Downloads Training Support Portal Partners Developers Community Search Sign In Language

**cloudera** Why Cloudera Products Services & Support Solutions Get Started

# QuickStart Downloads for CDH 5.5

Easy-to-deploy Apache Hadoop clusters for easy learning!

Cloudera QuickStart downloads contain complete Apache Hadoop clusters in the form of VMs or Docker images, including Cloudera Manager to manage them.

Get Started

QUICKSTART DOWNLOADS FOR CDH 5.5

SELECT A PLATFORM

www.nyc.gov/html/tlc/html/about/trip\_record\_data.shtml

NYC Resources 311 Office

**NYC** Taxi & Limousine Commission

Online Transactions (LARS) Printer Friendly Newsletter Sign-up Translate This Page Text Size

Home

About TLC

- » [TLC Mission Statement](#)
- » [Commission Room](#)
- » [TLC Facilities](#)
- » [Research and Statistics](#)
- » [Annual Reports](#)
- » [Employment Opportunities](#)
- » [Interagency MOUs](#)

TLC Rules and Local Laws

Licensing/Industry Information


Passenger Information

Frequently Asked Questions

TLC News

TLC Site Map

### TLC Trip Record Data



This dataset includes trip records from all trips completed in yellow and green taxis in NYC in 2014 and select months of 2015. Records include fields capturing pick-up and drop-off dates/times, pick-up and drop-off locations, trip distances, itemized fares, rate types, payment types, and driver-reported passenger counts. The data used in the

### Taxi News

The TLC and Fair Safe Streets pres  
Zero safety video  
Your Family Lives

< 2/5



trips10 - [Share] - Share - [~/Share]

trips10.csv

id	rate_code	store_and_fwd_flag	pickup_datetime	dropoff_datetime	passenger_count	trip_time_in_secs	trip_distance	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude
1	N		2013-01-01 15:11:48	2013-01-01 15:18:10	4	382	1.00	-73.978165	40.757977	-73.989838	40.751171
1	N		2013-01-06 00:18:35	2013-01-06 00:22:54	1	259	1.50	-74.006683	40.731781	-73.994499	40.75066
1	N		2013-01-05 18:49:41	2013-01-05 18:54:23	1	282	1.10	-74.004707	40.73777	-74.009834	40.726002
1	N		2013-01-07 23:54:15	2013-01-07 23:58:20	2	244	.70	-73.974602	40.759945	-73.984734	40.759388
1	N		2013-01-07 23:25:03	2013-01-07 23:34:24	1	560	2.10	-73.97625	40.748528	-74.002586	40.747868
1	N		2013-01-07 15:27:48	2013-01-07 15:38:37	1	648	1.70	-73.966743	40.764252	-73.983322	40.743763
1	N		2013-01-08 11:01:15	2013-01-08 11:08:14	1	418	.80	-73.995804	40.743977	-74.007416	40.744343
1	N		2013-01-07 12:39:18	2013-01-07 13:10:56	3	1898	10.70	-73.989937	40.756775	-73.86525	40.77063
1	N		2013-01-07 18:15:47	2013-01-07 18:20:47	1	299	.80	-73.980072	40.743137	-73.982712	40.735336

Commander  
Ant Build  
Database  
IDEtalk  
Maven Projects

2: Favorites

Text Data

6: TODO SBT Console Terminal

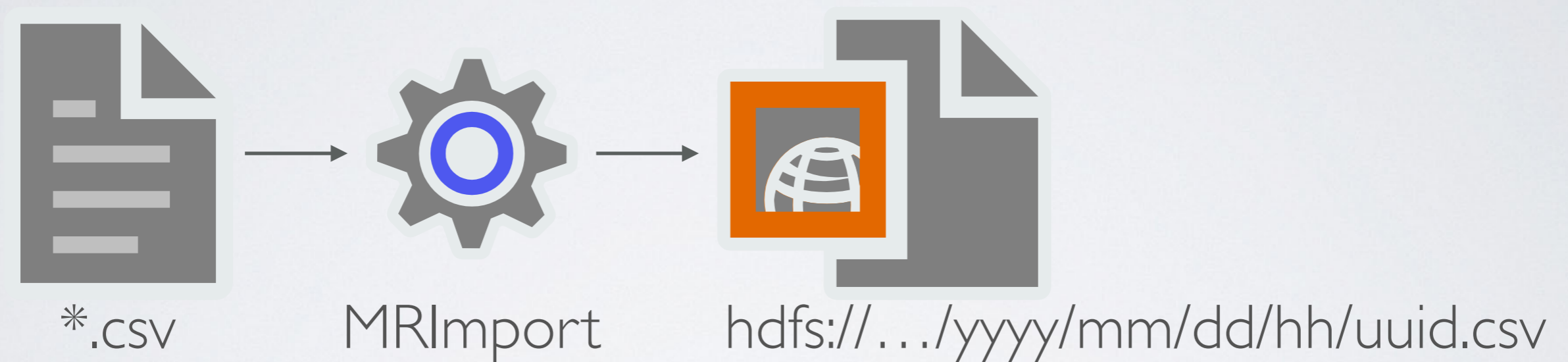
All files are up-to-date (28 minutes ago)

Event Log

n/a n/a 177 of 740M

~ 180 million entries - small :-)

# IMPORT BY TIME/SPACE





```
create external table if not exists trips (  
pickupdatetime string,  
dropoffdatetime string,  
pickupx double,  
pickupy double,  
dropoffx double,  
dropoffy double,  
passengercount int,  
triptime int,  
tripdist double,  
rc25 string,  
rc50 string,  
rc100 string,  
rc200 string  
) partitioned by (year int, month int, day int, hour int)  
row format delimited  
fields terminated by '\t'  
lines terminated by '\n'  
stored as textfile;
```



This repository Search

Explore Gist Blog Help



cloudera / **impyla**

Watch 23

Python client and Numba-based UDFs for Impala

152 commits

3 branches

2 releases

5 contributors



branch: master

**impyla** / +



Got sklearn back into working shape as well ...



laserson authored 2 days ago

latest commit 9df5e6f35e

bin	Fixed tab to space issues after large merges	22 days ago
examples	Fixed tab to space issues after large merges	22 days ago
impala	Got sklearn back into working shape as well	2 days ago
jenkins	Added Jenkins test script	5 days ago
thrift	Added support for HiveServer2 protocol V6	8 days ago



Desktop Help 10.0 - What x

help.arcgis.com/en/arcgisdesktop/10.0/help/index.html#/What\_is\_ArcPy/... tldr ☆

Apps BDMap t Y e 52 JS A E. V K GPU LingPipe

ArcGIS.com Esri.com

**ArcGIS Resource Center** Help Blogs Forums

Desktop 10

- Essential geoprocessing vocabulary
- Geoprocessing tools
- The geoprocessing framework
- Commonly used tools
- Finding tools
- Executing tools
- Managing tools and toolboxes
- Creating tools
- Sharing tools
- Geoprocessing with ModelBuilder
- Geoprocessing with Python
- Geoprocessing with ArcGIS Server
- The ArcPy site package
  - What is ArcPy?**
  - Essential ArcPy vocabulary
  - A quick tour of ArcPy
  - Functions
  - Classes
  - Mapping module
  - Geostatistical Analyst module
  - Spatial Analyst module
- Geoprocessing environment settings
- Geoprocessing tool reference
- Extensions
- ArcGIS Server

## What is ArcPy?

[Resource Center](#) » [Professional Library](#) » [Geoprocessing](#) » [The ArcPy site package](#)

ArcPy is a site package that builds on (and is a successor to) the successful arcgisscripting module. Its goal is to create the cornerstone for a useful and productive way to perform geographic data analysis, data conversion, data management, and map automation with Python.

This package provides a rich and native Python experience offering code completion (type a keyword and a dot to get a pop-up list of properties and methods supported by that keyword; select one to insert it) and reference documentation for each function, module, and class.

The additional power of using ArcPy within Python is the fact that Python is a general-purpose programming language. It is interpreted and dynamically typed and is suited for interactive work and quick prototyping of one-off programs known as scripts while being powerful enough to write large applications in. ArcGIS applications written with ArcPy benefit from the development of additional modules in numerous niches of Python by GIS professionals and programmers from many different disciplines.

### General Help

Python provides the facility of documentation strings. The functions and classes available in ArcPy use this method for the package documentation. One method for reading these messages and getting help is by using the command `help` provided by Python. Running the command with an argument displays the calling signature and the documentation string of the object.

```
>>> import arcpy
>>> help(arcpy)
```

Another method for getting help is the code completion provided by ArcPy. Anytime you

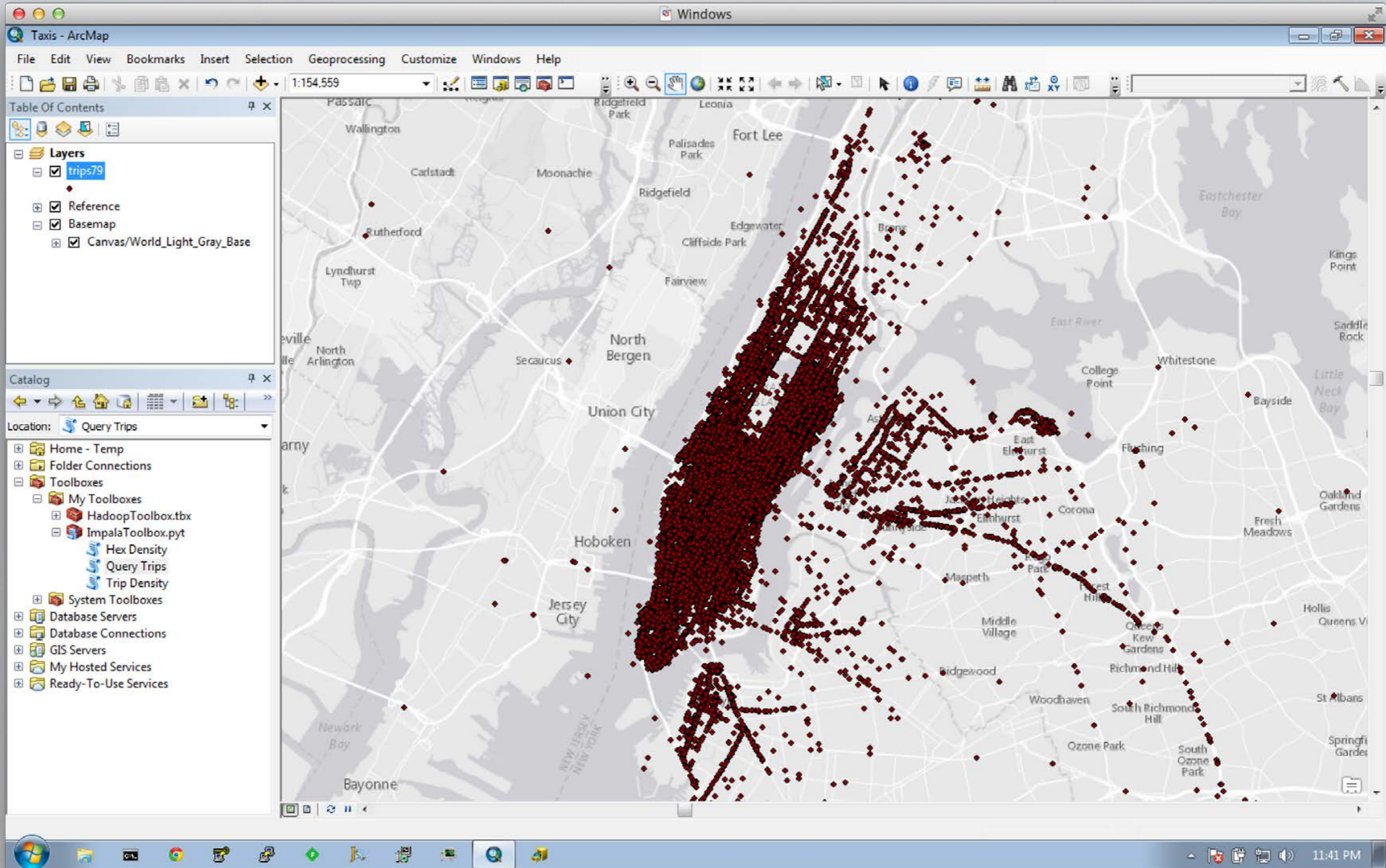


```

17
18 class TripTool(BaseTool):
19     def __init__(self):
20         super(TripTool, self).__init__()
21         self.label = "Query Trips"
22         self.description = "Tool to query trips table using Impala"
23         self.canRunInBackground = False
24
25     def getParameterInfo(self):
26         paramWhere = self.getParamString(name="in_where", displayName="Where", value="hour between 7 and 9")
27         return [paramWhere, self.getParamName(value="trips79"), self.getParamFC()]
28
29     def execute(self, parameters, messages):
30         try:
31             name = parameters[1].value
32             fc = "in_memory/" + name
33             self.deleteFC(fc)
34             spref = arcpy.SpatialReference(102100)
35             arcpy.management.CreateFeatureclass("in_memory", name, "POINT", spatial_reference=spref)
36             arcpy.management.AddField(fc, "PICKUP_DT", "TEXT")
37             arcpy.management.AddField(fc, "PASS_COUNT", "SHORT")
38             arcpy.management.AddField(fc, "TRIP_TIME", "SHORT")
39             arcpy.management.AddField(fc, "TRIP_DIST", "FLOAT")
40             conn = impala.dbapi.connect(host='quickstart', port=21050)
41             rows = conn.cursor()
42             hql = """select
43                 pickupx,
44                 pickupy,
45                 pickupdatetime,
46                 passengercount,
47                 triptime,
48                 tripdist
49             from trips
50             where {w}""".format(w=parameters[0].value)
51             hql = re.sub(r'\s+', ' ', hql)
52             rows.execute(hql)
53             with arcpy.da.InsertCursor(fc,
54                                     ['SHAPE@XY', 'PICKUP_DT', 'PASS_COUNT', 'TRIP_TIME', 'TRIP_DIST']) as cursor:
55                 for row in rows:
56                     cursor.insertRow([(row[0], row[1]), row[2], row[3], row[4], row[5]])
57                 del row
58             del rows
59             parameters[2].value = fc
60         except:
61             arcpy.AddMessage(traceback.format_exc())

```





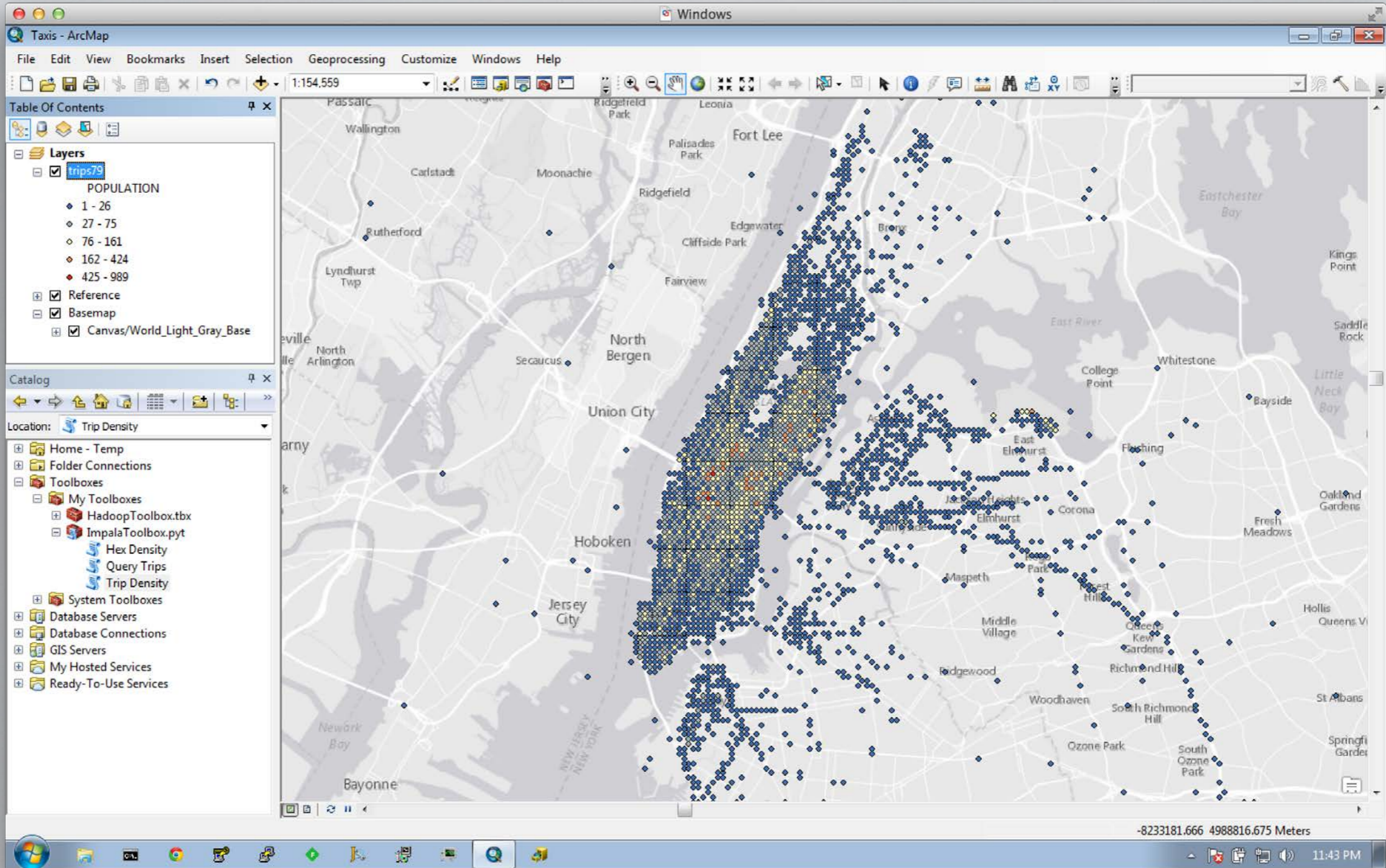


```

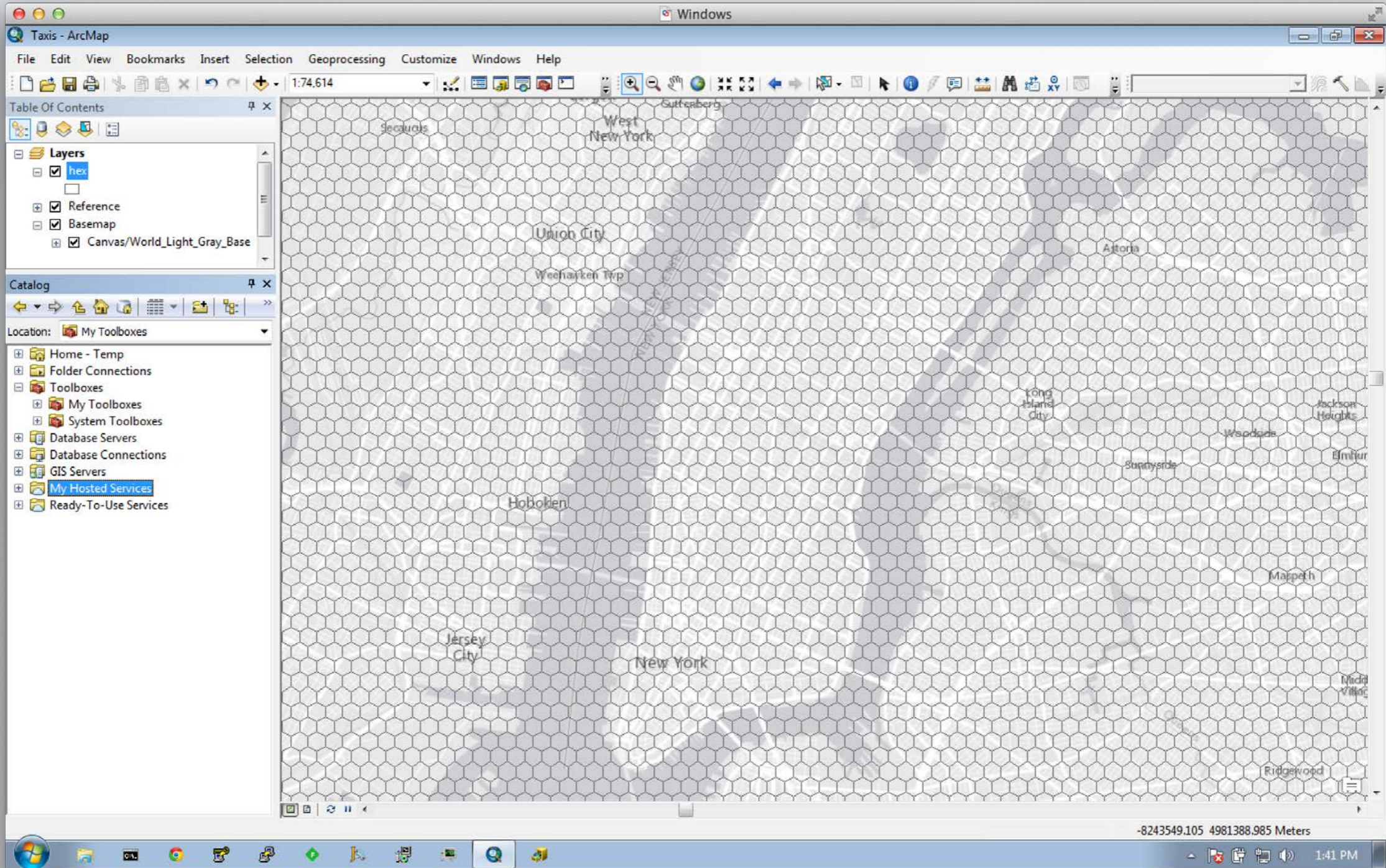
76 def execute(self, parameters, messages):
77     try:
78         cell1 = float(parameters[0].value)
79         cell2 = cell1 * 0.5
80         where = parameters[1].value
81         name = parameters[2].value
82
83         fc = "in_memory/" + name
84         self.deleteFC(fc)
85         spref = arcpy.SpatialReference(102100)
86         arcpy.management.CreateFeatureclass("in_memory", name, "POINT", spatial_reference=spref)
87         arcpy.management.AddField(fc, "POPULATION", "LONG")
88         conn = impala.dbapi.connect(host='quickstart', port=21050)
89         rows = conn.cursor()
90         hql = """select
91             T.X*{c1}+{c2} as X,
92             T.Y*{c1}+{c2} as Y,
93             count(*) AS POPULATION from (
94             select
95                 cast(floor(pickupx/{c1}) as int) as X,
96                 cast(floor(pickupy/{c1}) as int) as Y
97             from trips where {w}) T
98             group by T.X,T.Y
99             """.format(w=where, c1=cell1, c2=cell2)
100         hql = re.sub(r'\s+', ' ', hql)
101         rows.execute(hql)
102         with arcpy.da.InsertCursor(fc, ['SHAPE@XY', 'POPULATION']) as cursor:
103             for row in rows:
104                 cursor.insertRow([(row[0], row[1]), row[2]])
105             del row
106         del rows
107         parameters[3].value = fc
108     except:
109         arcpy.AddMessage(traceback.format_exc())
110

```





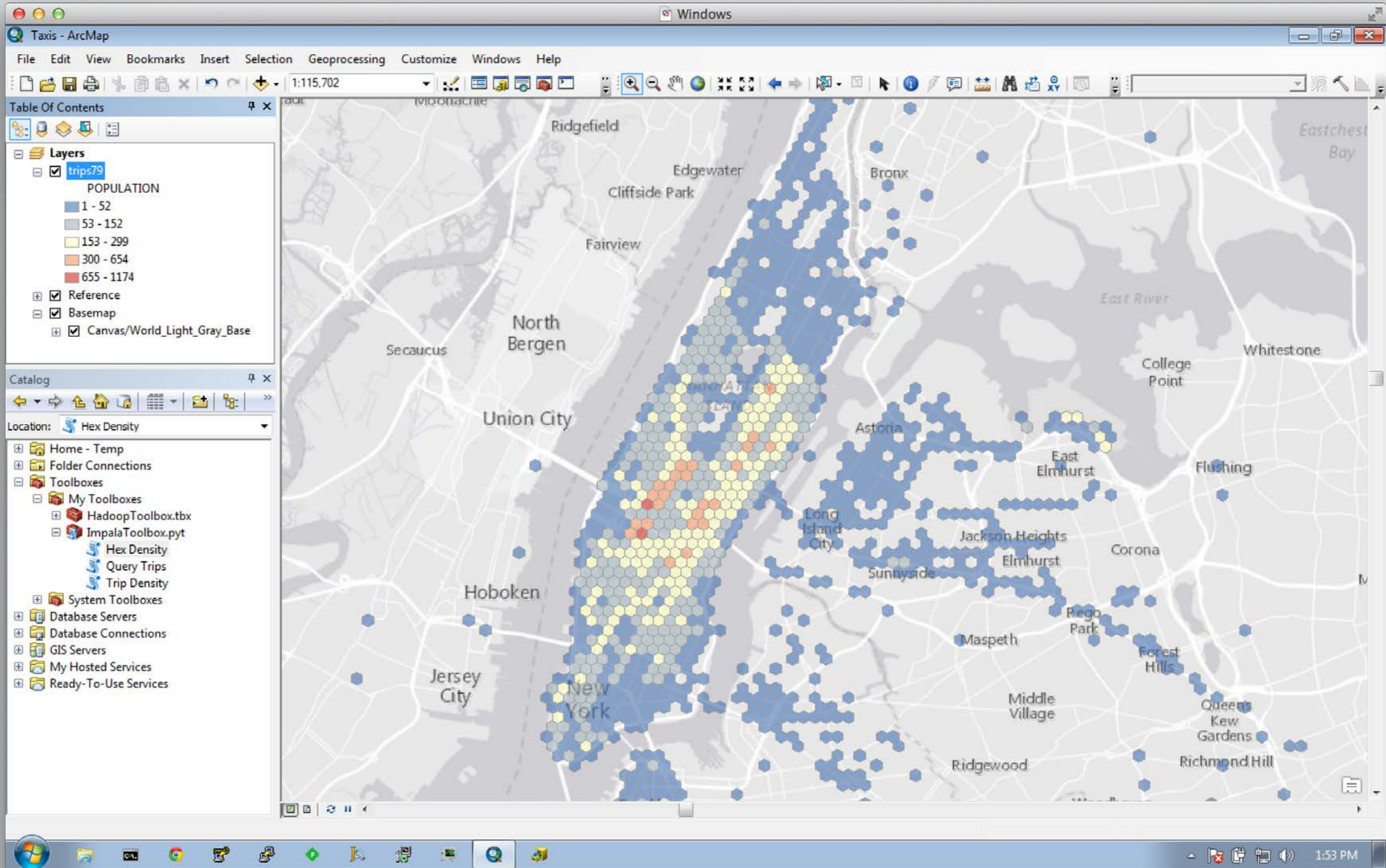






```
cloudera-quickstart-vm-5.1.0-1-vmware
Applications Places System Home - Cloudera Man... cloudera@quickstart:~ Wed Oct 15, 10:48 AM cloudera
cloudera@quickstart:~
File Edit View Search Terminal Help
Welcome to the Impala shell. Press TAB twice to see a list of available commands.
Copyright (c) 2012 Cloudera, Inc. All rights reserved.
(Shell build version: Impala Shell v1.4.0-cdh5-INTERNAL (e801bd8) built on Sat Jul 12 06:45:04 PDT 2014)
Query: invalidate metadata
Returned 0 row(s) in 3.75s
[quickstart.cloudera:21000] > describe trips;
Query: describe trips
+-----+
| name      | type  | comment |
+-----+
| pickupdatetime | string |         |
| dropoffdatetime | string |         |
| pickupx      | double |         |
| pickupy      | double |         |
| dropoffx     | double |         |
| dropoffy     | double |         |
| passengercount | int    |         |
| triptime     | int    |         |
| tripdist     | double |         |
| rc25         | string |         |
| rc50         | string |         |
| rc100        | string |         |
| rc200        | string |         |
| year         | int    |         |
| month        | int    |         |
| day          | int    |         |
| hour         | int    |         |
+-----+
Returned 17 row(s) in 14.18s
[quickstart.cloudera:21000] > select pickupx,pickupy,rc25,rc100 from trips limit 5;
Query: select pickupx,pickupy,rc25,rc100 from trips limit 5
+-----+-----+-----+-----+
| pickupx      | pickupy      | rc25      | rc100      |
+-----+-----+-----+-----+
| -8235510.550453553 | 4977585.413571722 | 4736|1489 | 1184|372 |
| -8231890.885890919 | 4975433.207053881 | 4678|1573 | 1170|393 |
| -8236461.775502382 | 4976817.795186556 | 4715|1467 | 1179|366 |
| -8214370.645193438 | 4960174.878029009 | 4271|1977 | 1068|494 |
| -8234755.581666993 | 4977389.204397635 | 4730|1507 | 1183|376 |
+-----+-----+-----+-----+
Returned 5 row(s) in 1.06s
[quickstart.cloudera:21000] > █
```







# STATISTICAL SIGNIFICANCE ?

# HOTSPOT ANALYSIS

The Getis-Ord local statistic is given as:

$$G_i^* = \frac{\sum_{j=1}^n w_{i,j} x_j - \bar{X} \sum_{j=1}^n w_{i,j}}{S \sqrt{\frac{n \sum_{j=1}^n w_{i,j}^2 - \left( \sum_{j=1}^n w_{i,j} \right)^2}{n-1}}} \quad (1)$$

where  $x_j$  is the attribute value for feature  $j$ ,  $w_{i,j}$  is the spatial weight between feature  $i$  and  $j$ ,  $n$  is equal to the total number of features and:

$$\bar{X} = \frac{\sum_{j=1}^n x_j}{n} \quad (2)$$

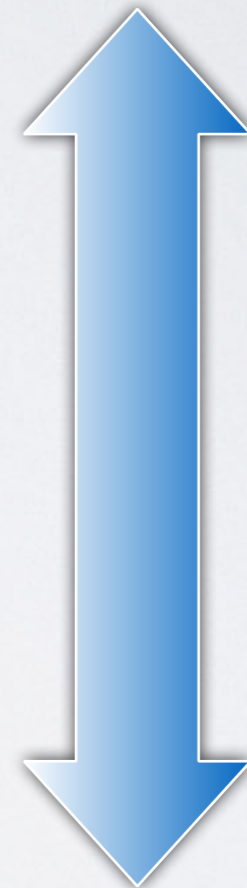
$$S = \sqrt{\frac{\sum_{j=1}^n x_j^2}{n} - (\bar{X})^2} \quad (3)$$

The  $G_i^*$  statistic is a  $z$ -score so no further calculations are required.



# PROCESSING EVOLUTION

- ✘ Transaction - Batch
- ✘ Operational - Dashboard
- ✘ Analytics - Exploration
- ✘ Intelligent - Realtime / Predictive



# WHAT IS NEXT ?

- In Memory
- Native Spatial Index In NoSQL DB
- Native Spatial Types (Point, Line, ...)
- Out-of-the-box Spatial Operators / Operations
- Distributed/Disconnected GPU Integration
- Visualization via Gamification



# Q&A

**Mansour Raad**  
**[thunderheadxplorer.blogspot.com](http://thunderheadxplorer.blogspot.com)**  
**[mraad@esri.com](mailto:mraad@esri.com)**  
**@mraad**

