



Scaling Spatio-Temporal Analytics with GeoMesa



Maryland Data Science Meetup -- August 18, 2014
Andrew Hulbert -- ahulbert@ccri.com

Agenda

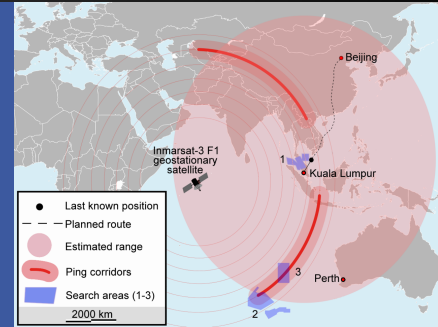
- Spatio-Temporal Data
- Distributed Indexing Strategies
- GeoMesa Architecture
- SpatioTemporal Analytics
- Roadmap & Future Work

What is Spatio-Temporal Data?

- Satellite Imagery
- FAA Flight Information
- Twitter & Social Media
- GPS-Enabled Apps
- Network Traffic & Clickstreams



twitter



What is GeoMesa?

- **Distributed** Spatio-Temporal Database
 - Apache Accumulo
 - Apache Hadoop
- **Standards-based** Data Access
 - GeoTools
 - GeoServer
- LocationTech **Open Source**
 - Eclipse Foundation

Open Source Geo



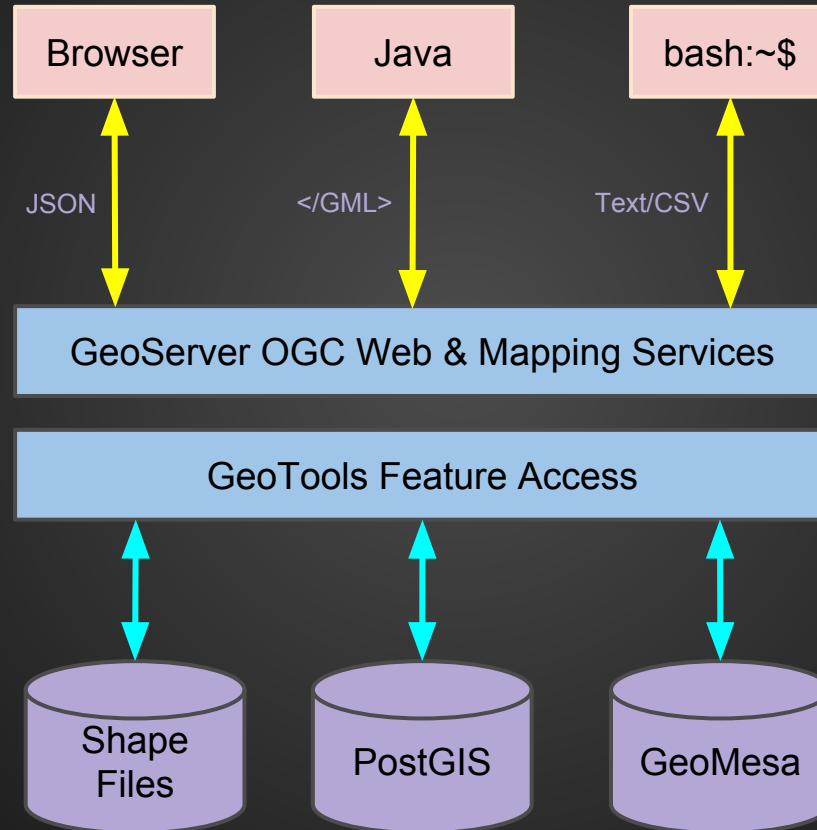
Simple Features

```
{
  "id" : "c11b0bdc-a400-42b3-9177-81f835bff6f5",
  "geometry" : {
    "coordinates" : [-74.189, 40.7728],
    "type" : "Point"
  },
  "type" : "Feature",
  "properties" : {
    "dtg" : "2014-07-08T04:08:43.000+0000",
    "text" : "\"This is our city!\" -David Ortiz",
    "tweet_id" : 010101001010101001010101001,
    "user_id" : 010101001,
    "user_name" : "Fresh Prince of Belair"
  }
}
```

Common Query Language (CQL)

```
BBOX(geom, 40, 30, 45, 35) AND
name = "Person 1" AND
age = 30 AND
hobby <> "texting"
```

Tiered Client Architecture



41.26289, 49.03574

Database Evolution

- PostGIS
- R-Trees
- Vertical Scaling



- Hadoop
- Key Value Stores
- Linear Scaling



Software Engineering 101

Hop on over from PostGIS....

```
val params = Map("dbtype"    -> "postgis",  
                 "host"      -> "localhost",  
                 "port"      -> 5432,  
                 "schema"    -> "public",  
                 "database"  -> "gdelt",  
                 "user"      -> "postgres",  
                 "passwd"    -> "postgres")  
val datastore = DataStoreFinder.getDataStore(params)
```

to GeoMesa

```
val params = Map("instanceId" -> "accumulo",  
                 "zookeepers" -> "zoo1:2181",  
                 "user"       -> "geomesa",  
                 "password"   -> "supersecret",  
                 "auths"      -> "",  
                 "tableName"  -> "gdelt")  
val datastore = DataStoreFinder.getDataStore(params)
```

Accumulo: Distributed Database

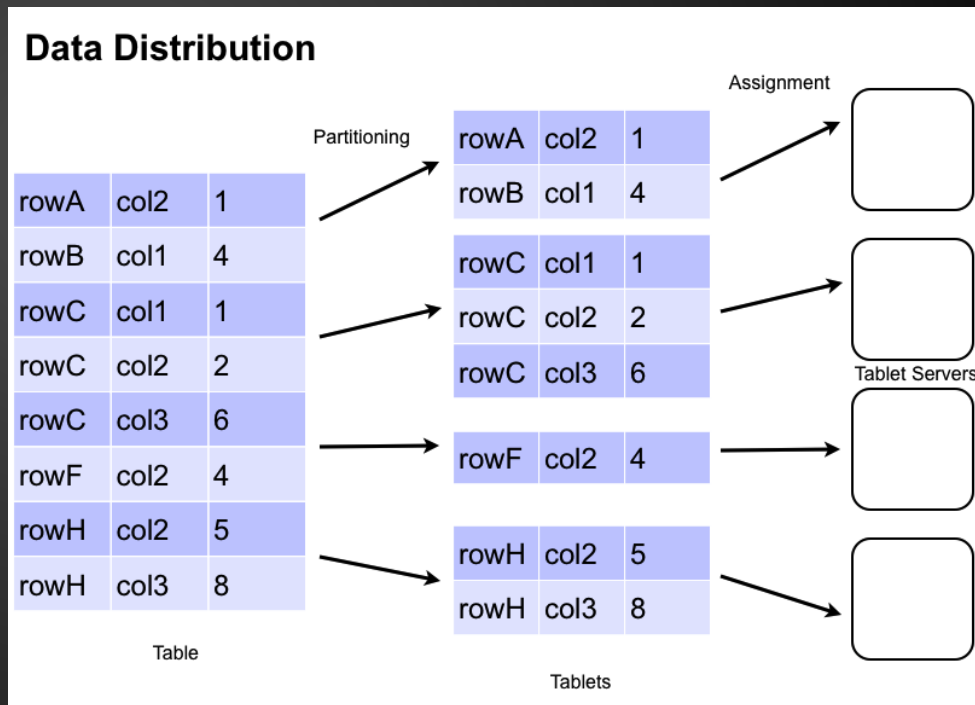
- Distributed Key/Value Store
 - 5-tuple Key
- Lexicographic Ordering
- Schema-less...



Key					Value
Row ID	Column			Timestamp	
	Family	Qualifier	Visibility		

Apache Accumulo

- Tables composed of Tablets
- Tablet Servers
- Partitioning Data
- Parallel Queries



(E)CQL in Accumulo?

- BBOX

- `bbox(geom, -74.1, 40.7, -74.2, 40.8)`

- INTERSECTS

- `intersects(geom, Polygon((-74.1 40.7, -74.2 40.7, -74.2 40.8, -74.1 40.8, -74.1 40.7)))`

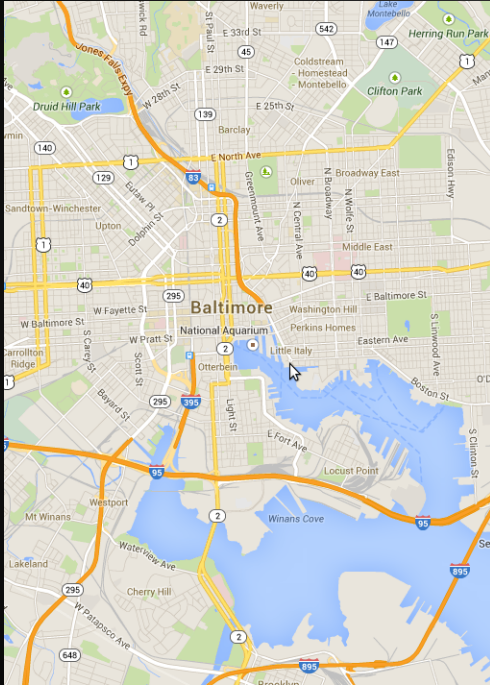
- DWITHIN

- `dwithin(geom, Point(-74.00 40.71), 5000, meters)`

- Temporal Predicates

- `dtg between 2014-01-02 AND 2014-01-06`

How will we do this?

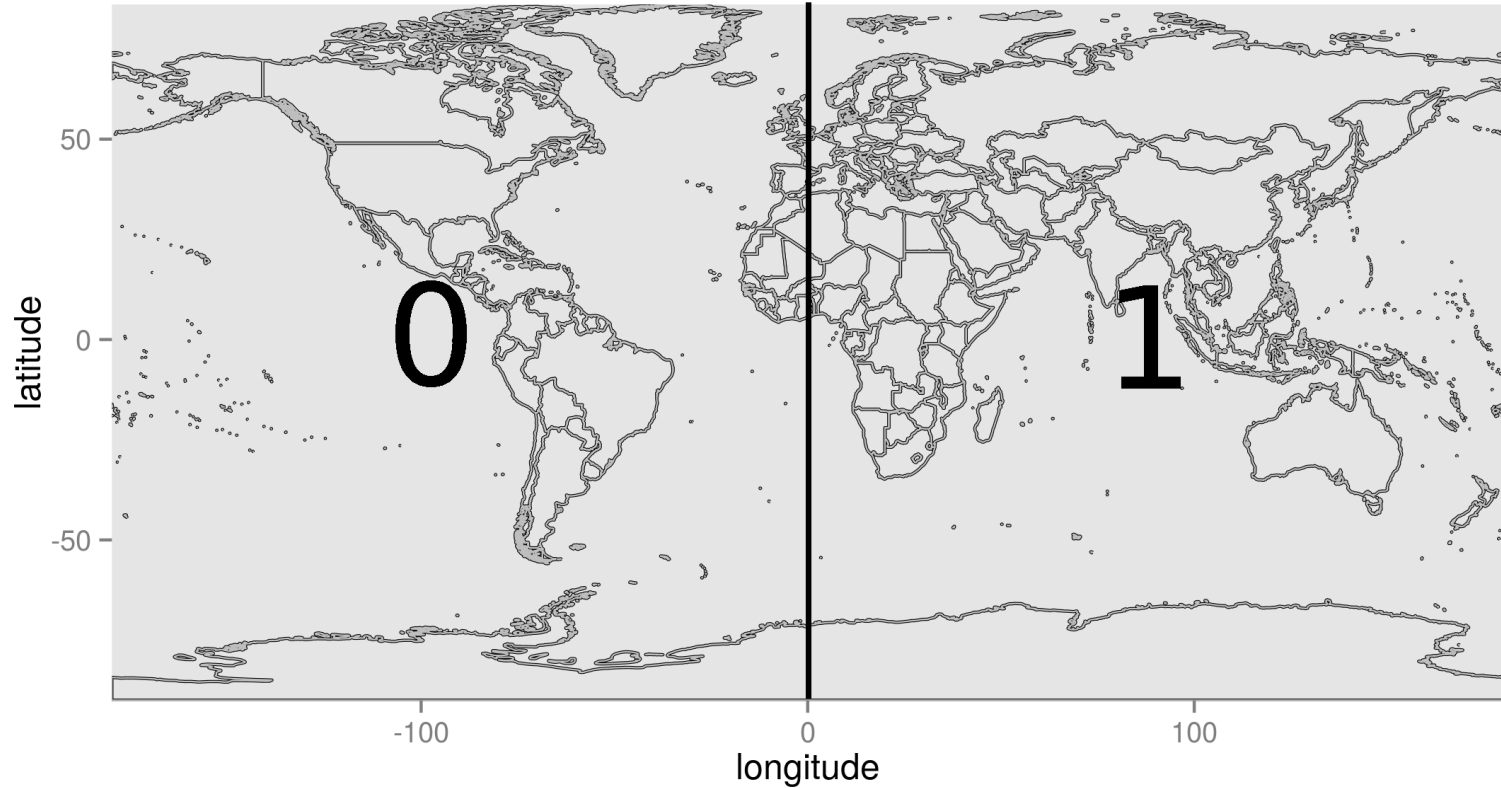


Key					Value
Row ID	Column			Timestamp	
	Family	Qualifier	Vis		

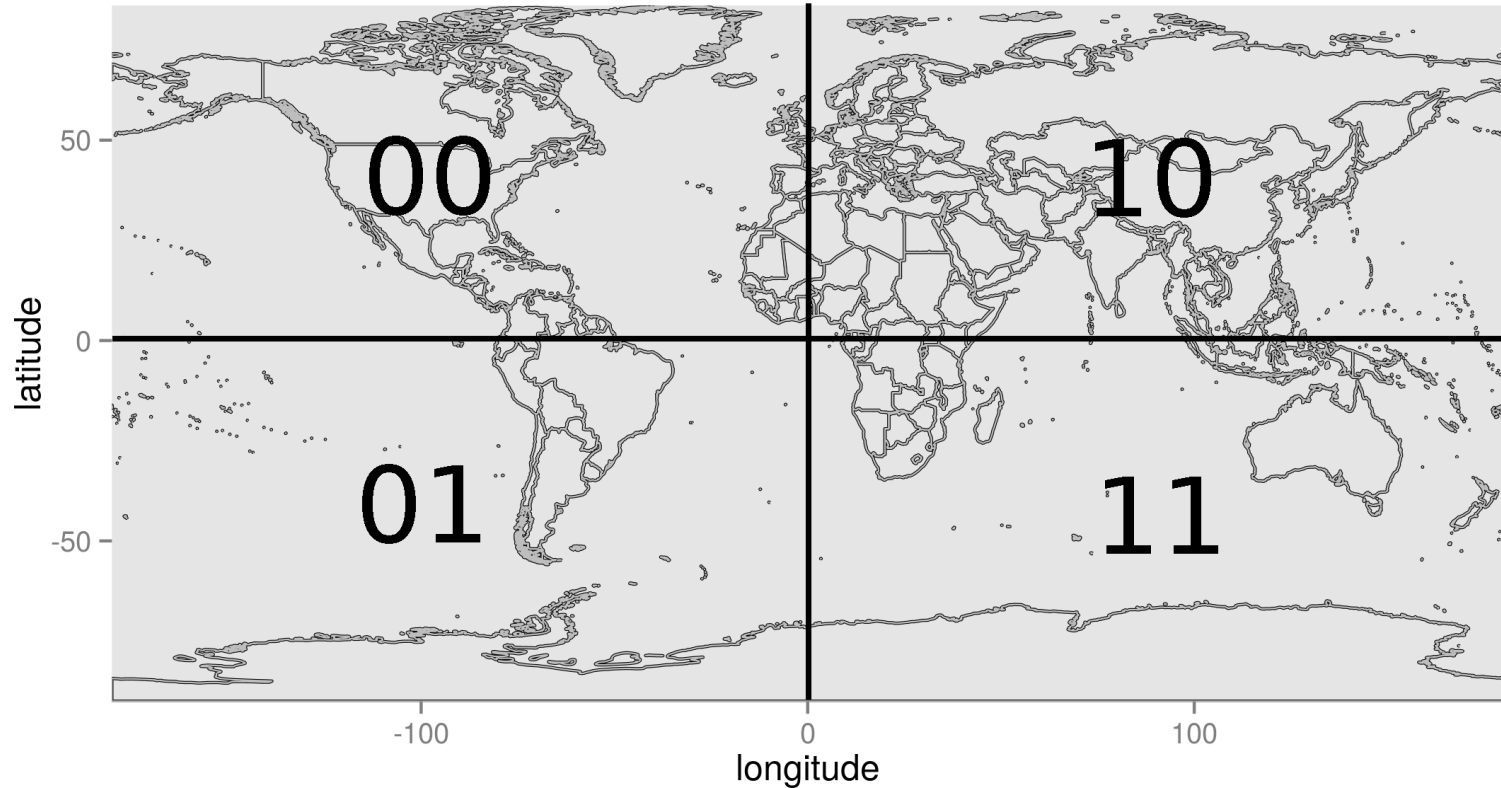
Space Filling Curves

- Z-order Curve / GeoHash Algorithm
 - Encode multiple dimensions
- Interleaved Lon-Lat bits
 - Left to Right binary encoding
 - Alternate Splitting Rectangles
 - Left -> 0
 - Right -> 1
 - Top -> 0
 - Bottom -> 1

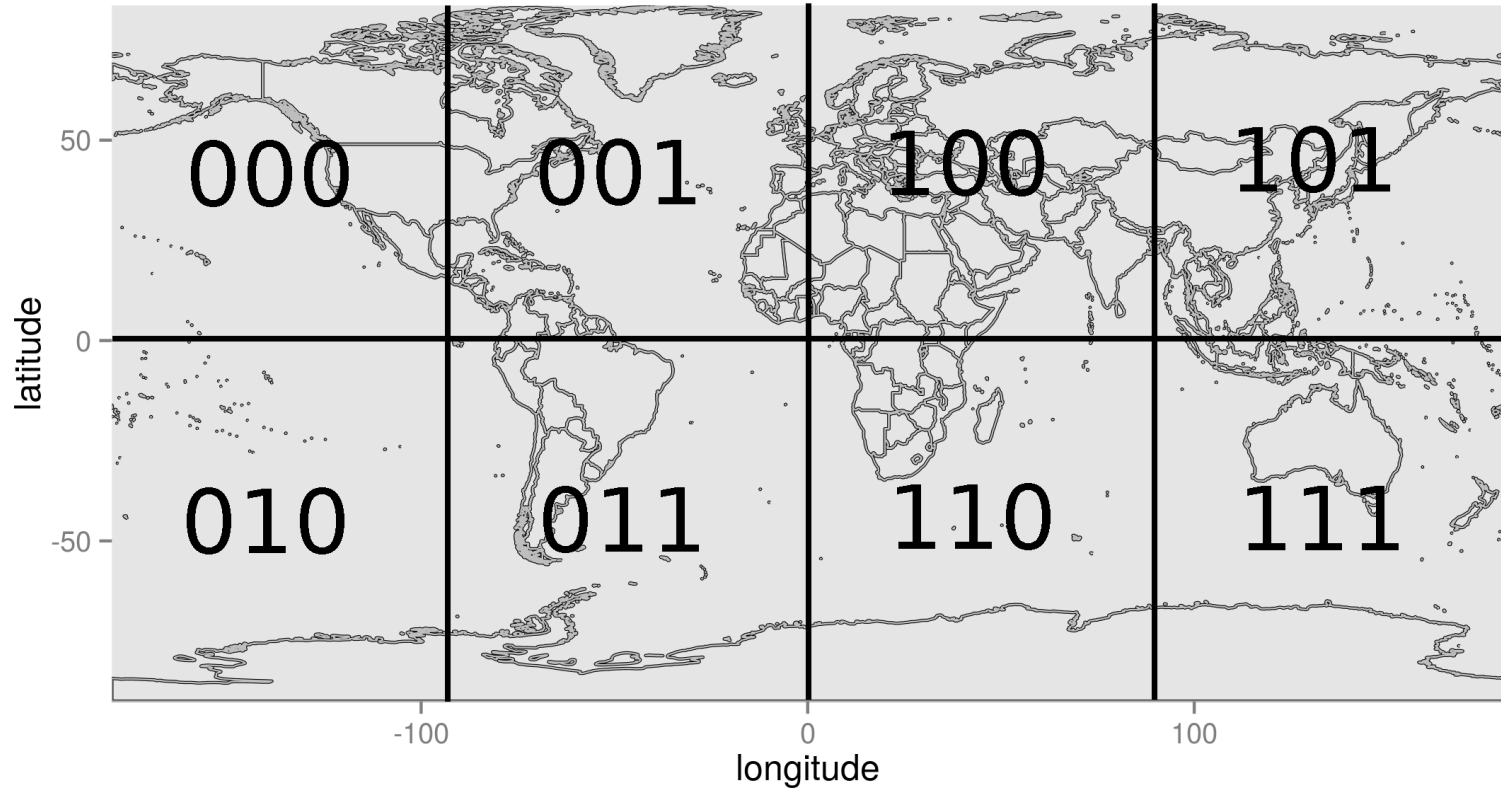
Linearize the Keyspace



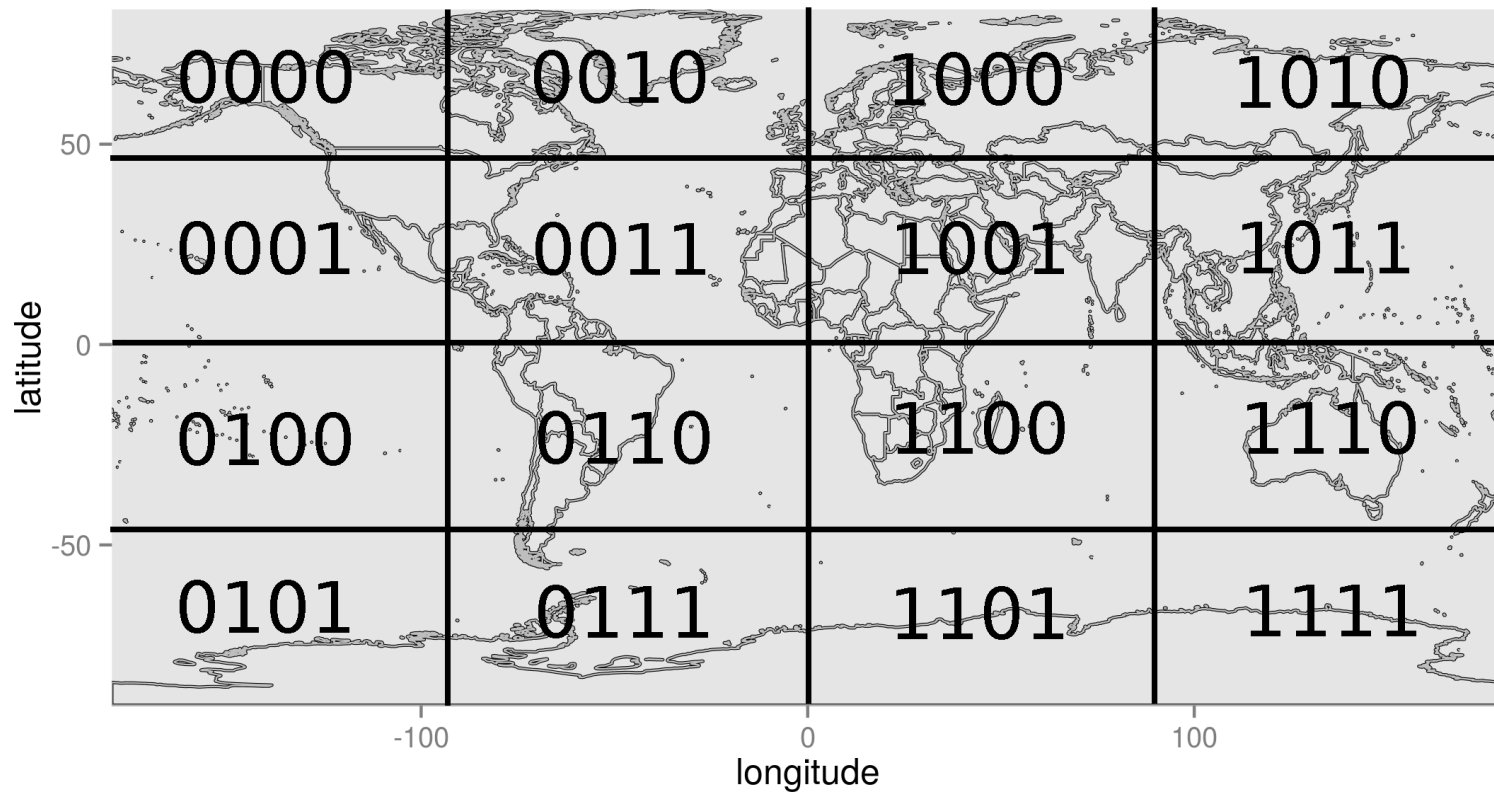
Linearize the Keyspace



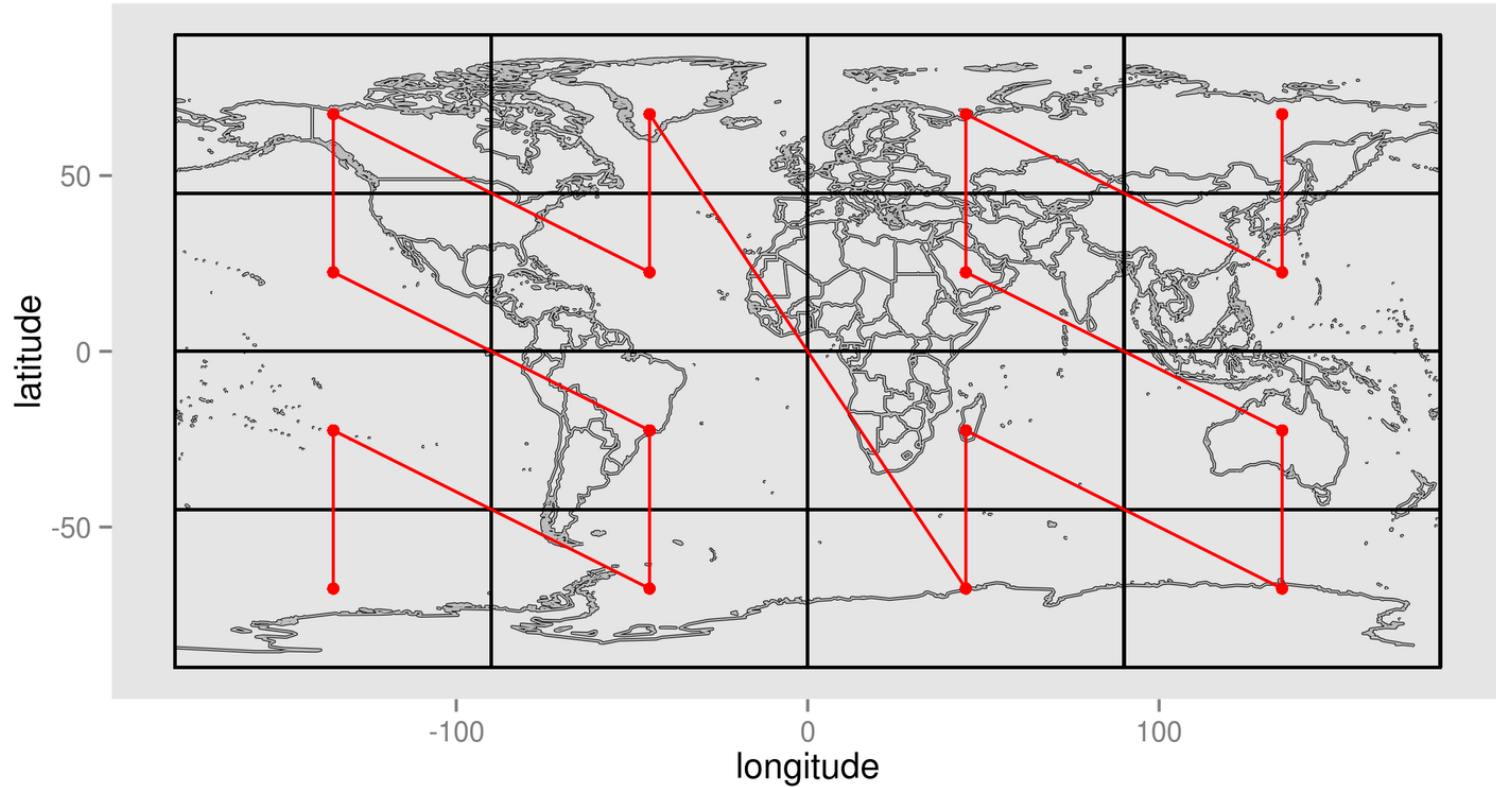
Linearize the Keyspace



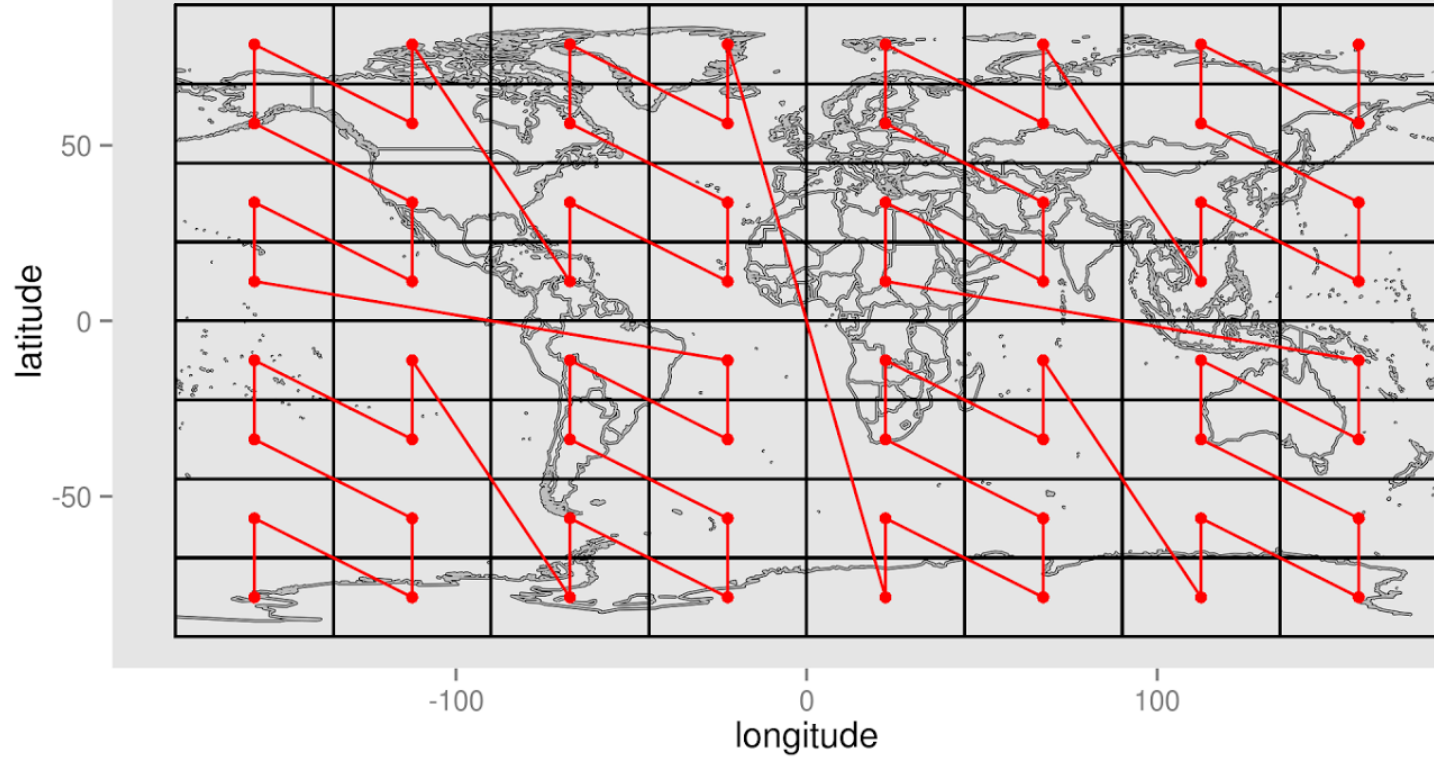
Linearize the Keyspace



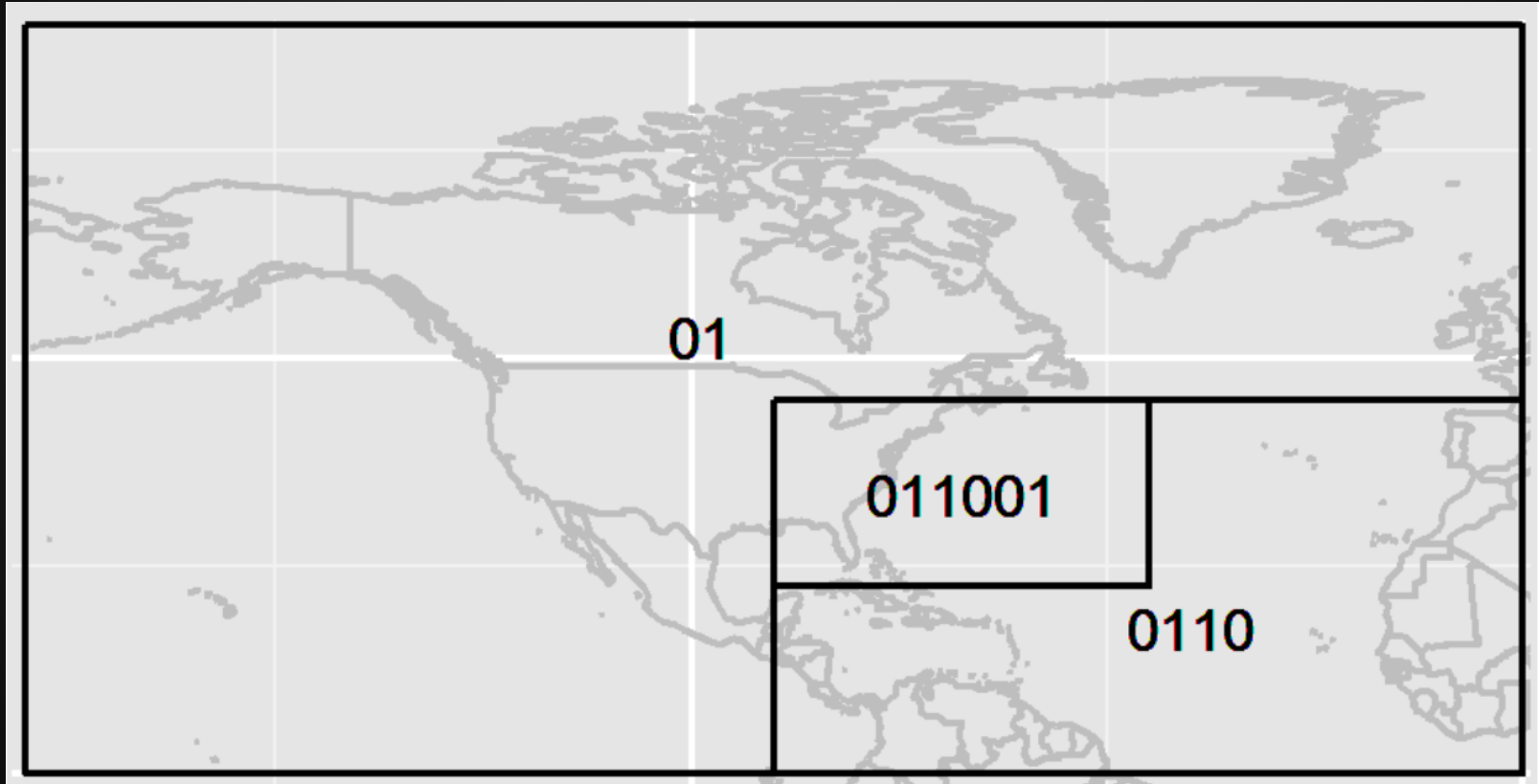
Linearize the Keyspace



Higher Precision



Locality and Containment

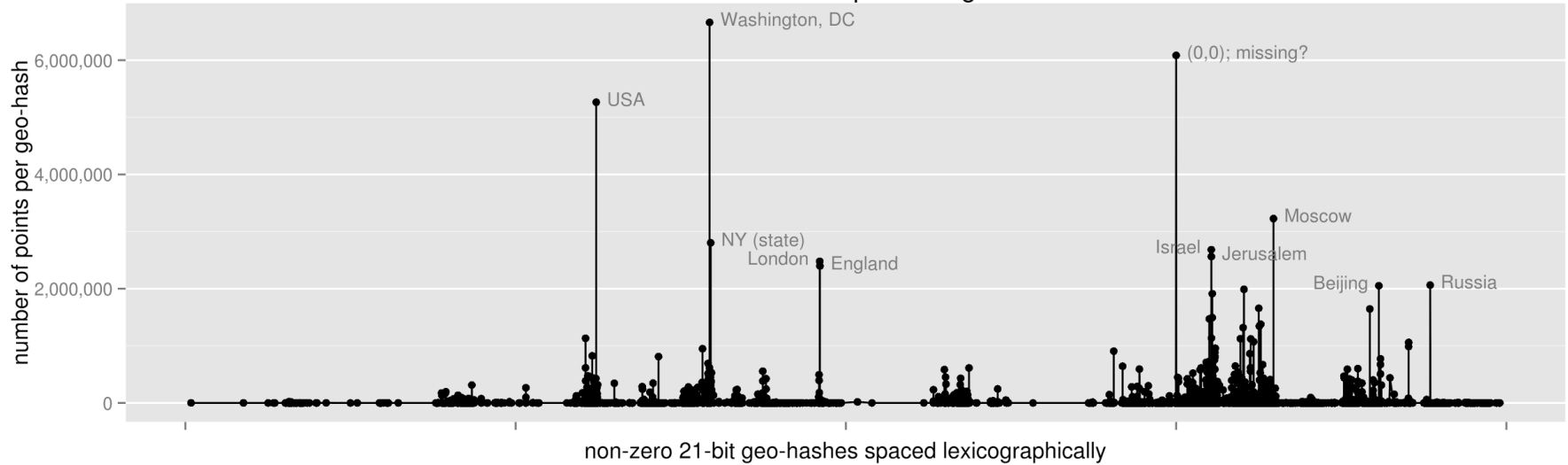


Locality and Containment

Baltimore (downtown)	-76.6167, 39.2833	dqcx2x33v
Fells Point	-76.5928, 39.2831	dqcx3p6bd
Federal Hill	-76.6100, 39.2789	dqcx2ww21
Dundalk	-76.5053, 39.2658	dqcx7h4ec
Columbia	-76.8569, 39.2036	dqcr50f98
Annapolis	-76.5012, 38.9729	dqctenv98
Philadelphia	-75.1667, 39.9500	dr4e389m8
New York	-74.0059, 40.7127	dr5regw89

Geospatial Hotspots

distribution of GDELT points to geo-hashes



San Fran



Sharding San Fran



Key Structure

Partition ID	L1 Spatial Hash	L1Temporal Hash	L2 Spatial Hash	L2 Temporal Hash	ID	Geometry	Time
--------------	-----------------	-----------------	-----------------	------------------	----	----------	------


Index	Row Key			Column Family		Column Qualifier	Value	
	001	9q8yy	201010	u	10	tweet.27963	POINT (-151,37)	20101009T23:47:07.234
	001	9q8yy	201010	v	15	tweet.12000	POINT (-151,37)	20101009T23:47:07.234
	001	gq89z	201112	x	25	tweet.32773	POINT (-151,37)	20101009T23:47:07.2343

	001	t8vtw	200805	a	10	tweet.65432	POINT (-151,37)	20101009T23:47:07.234
Data	Row Key			Column Family		Column Qualifier	Value	
	001	9q8yy	201010	tweet.27963		#OGC, #BigData	id=xyz green 48	
	001	9q8yy	201010	tweet.12000		#Accumulo	id=abc yellow 12	
	001	gq89z	201112	tweet.32773		#BigTable	state=red	
	
	001	t8vtw	200805	tweet.65432		#Geohash	value=conference	

Building an Architecture

- GeoServer Integration
- Security
- Command Line Tools
- Indexing, Querying Planning
- Compute + Analytics

GeoServer Integration

**GeoServer**

Logged in as admin. [Logout](#)

About & Status

- Server Status
- GeoServer Logs
- Contact Information
- About GeoServer

GeoMesa

- Data Stores
- Hadoop Status
- Configuration
- Actions

Data

- Layer Preview
- Workspaces
- Stores
- Layers
- Layer Groups
- Styles

Services

- WCS
- WFS
- WMS
- WPS

Settings

- Global
- JAI
- Coverage Access

Tile Caching

- Tile Layers
- Caching Defaults
- Gridsets
- Disk Quota

Security

- Settings
- Authentication
- Passwords
- Users, Groups, Roles
- Data
- Services

Demos

GeoMesa Data Stores

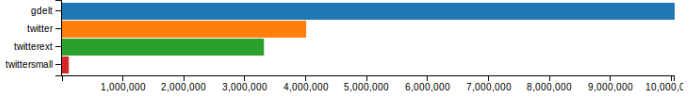
<< < 1 > >> Results 1 to 11 (out of 11 items)

datatype	workspace	name	type	enabled	
<input type="checkbox"/>		geomesa	geomesa	Accumulo Feature Data Store	<input checked="" type="checkbox"/>

<< < 1 > >> Results 1 to 11 (out of 11 items)

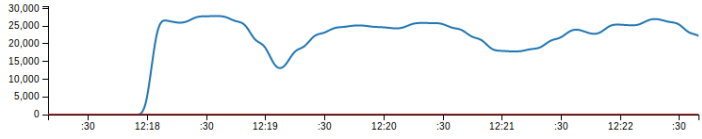
Data Visualization

Number of Entries



Store	Number of Entries
gdel	10,039,507
twitter	4,000,000
twitterext	3,500,000
twitterstmal	500,000

Ingest Rate (entries/s)



Time	Ingest Rate (entries/s)
12:18:00	0
12:18:30	25,000
12:19:00	25,000
12:19:30	15,000
12:20:00	25,000
12:20:30	25,000
12:21:00	15,000
12:21:30	25,000
12:22:00	25,000
12:22:30	25,000

Data Store - geomesa:geomesa

Feature: gdel

Table	Number of Tablets	Number of Splits	Total Entries	Total Size (MB)
Record Table	10	30	10,039,507	450.6
GeoSpatial Index	22	96	12,095,002	501.83
Attribute Index	0	0	0	0

Bounds: -180.0:180.0, -90.0:90.0

[Feature Attributes](#)

Command Line Tools

- Feature Ingest
- Query & Export
 - Shape Files, csv, tsv
- Data Management
 - List, Describe, Add, Remove
- Explain Queries

Query Planning and Indexing

- SpatioTemporal Index
- Secondary Index on Attributes
- Queries:
 - How do presidents' word choices change based on the state they're visiting?
 - Where are people eating in Baltimore before Orioles games?

SpatioTemporal Index

- Political Language

```
bbox("United States") AND  
time between 2000-2014
```



- O's Game

```
dwithin("Camden Yards", 3, "km") AND  
text like %food% AND  
(time between t1 AND t2 OR  
time between t3 and t4, ...)
```



Attribute Index

- Political Language

`bbox("United States") AND`
`time between 2000-2014`
`user_id in (1, 2, 3, 4, 5, etc)`



- O's Game

- What are all the user id's?
- Full table scan



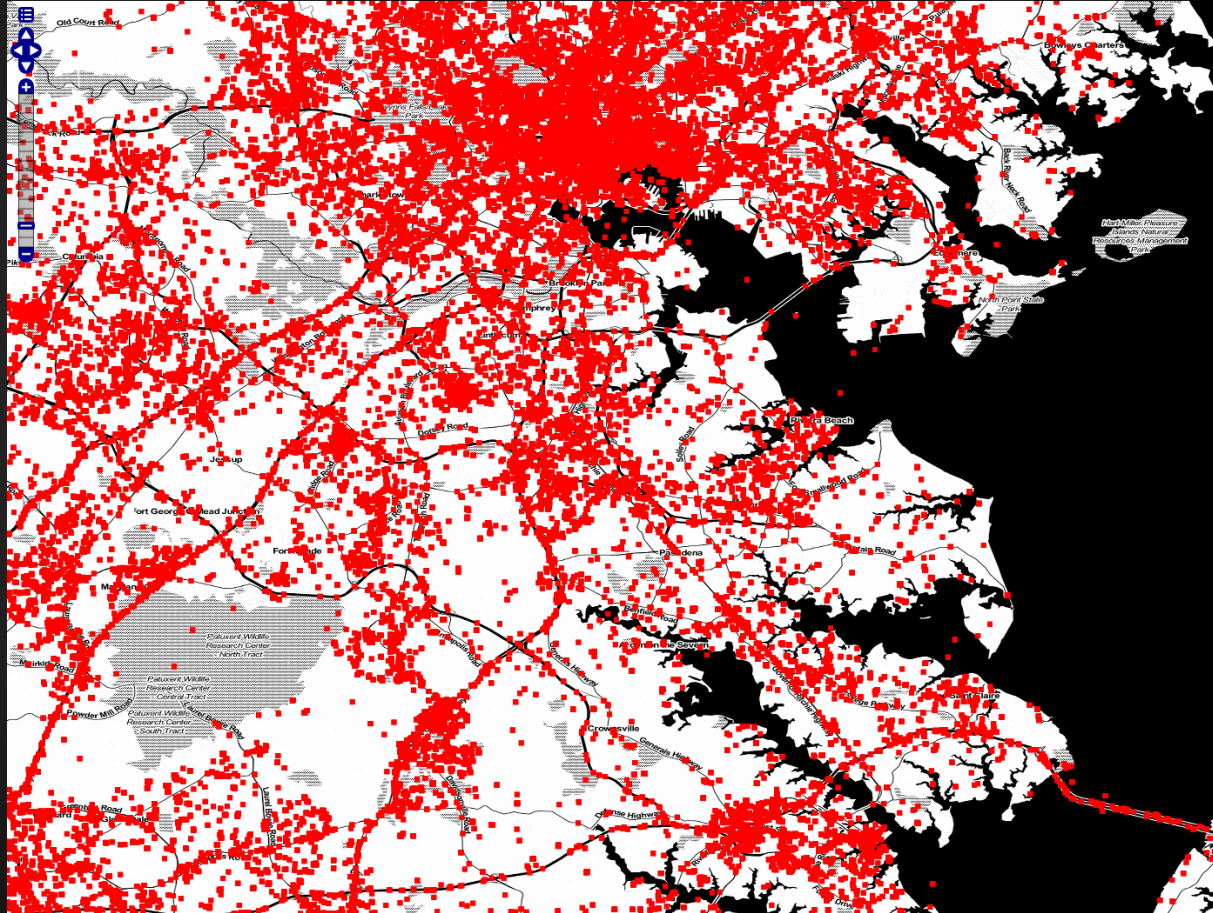
Query Planning

- Table Statistics
- Query Statistics
- Attribute Cardinality
- Iterators
 - SpatioTemporal Intersecting
 - Attribute Filtering
 - ST Index Only
 - Attribute Index

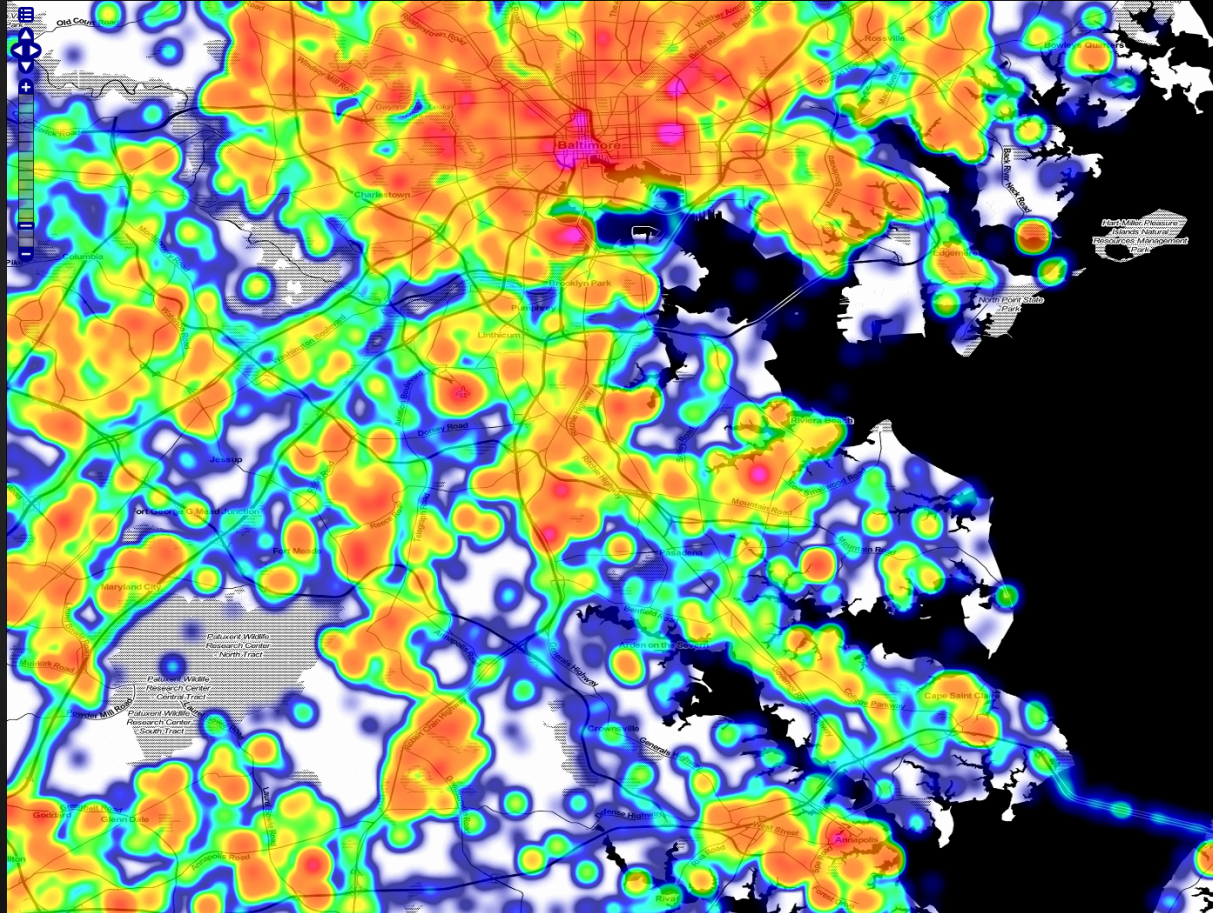
Compute & Analytics

- Queries & Mapping
- Accumulo Iterators
- GeoServer WPS Processes
 - AaaS (Analytics as a Service)
- Scalding
- Spark
- Data Export

CQL Queries



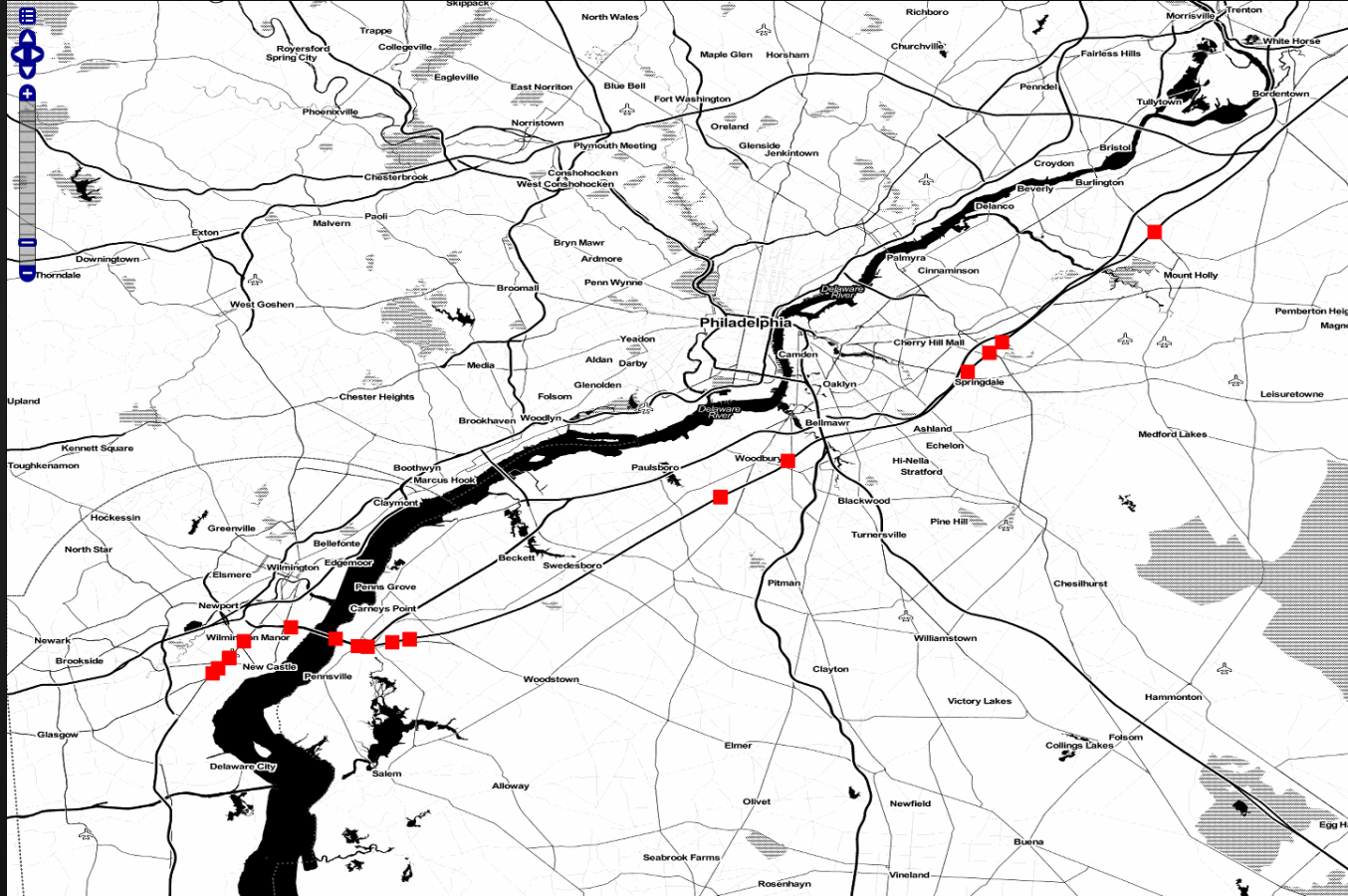
Parallel Density Calculations



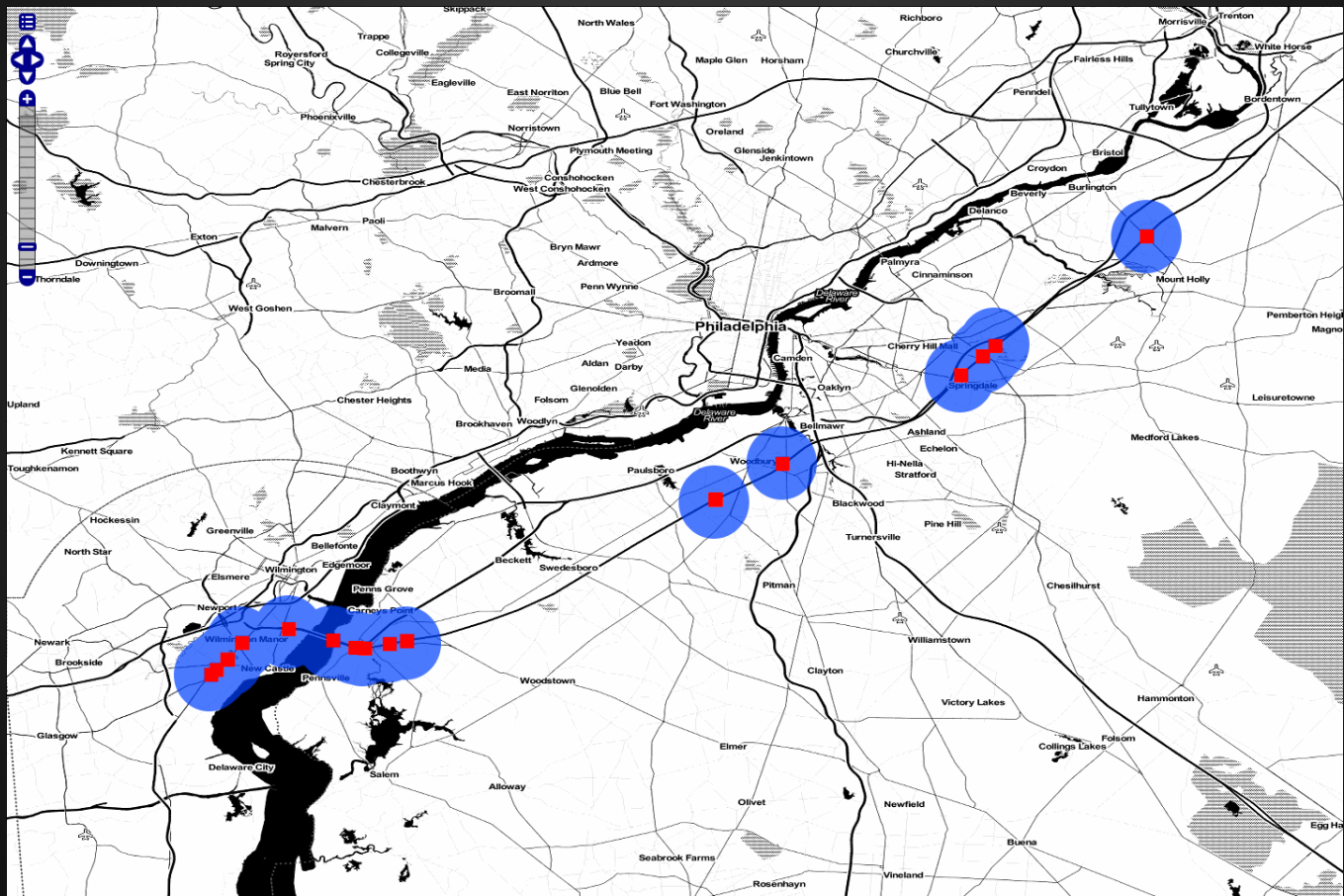
Time-Interpolated Queries

- Spatio-temporal gradient
 - Interpolate space (max speed, max time)
 - Interpolate time
- Co-Travelers
- Alias detection
 - Clickstream by phone IP addr
- Candidate Interactions
 - Subsequent scoring
 - Motion Models

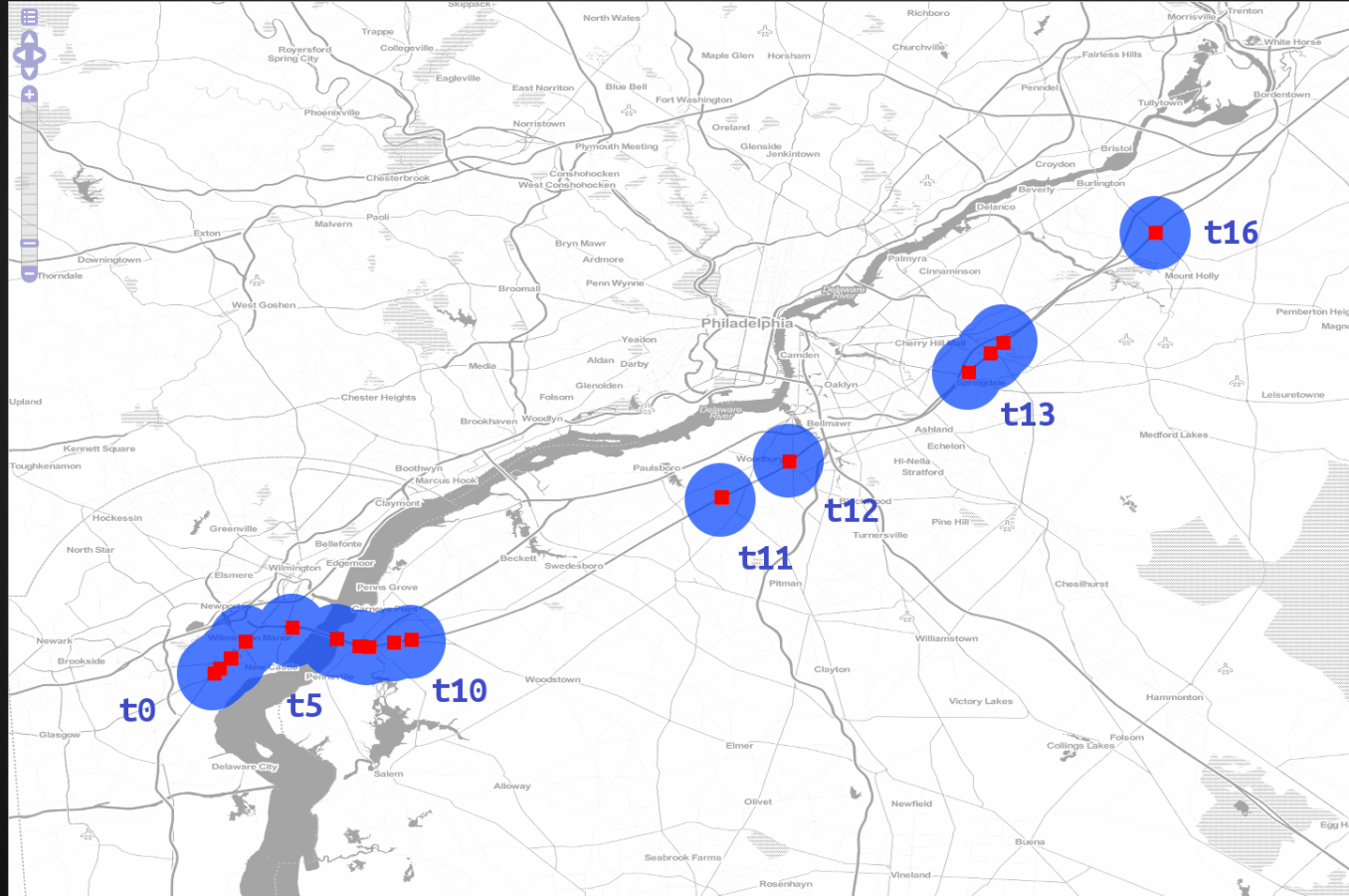
Tweeting the NJ Turnpike



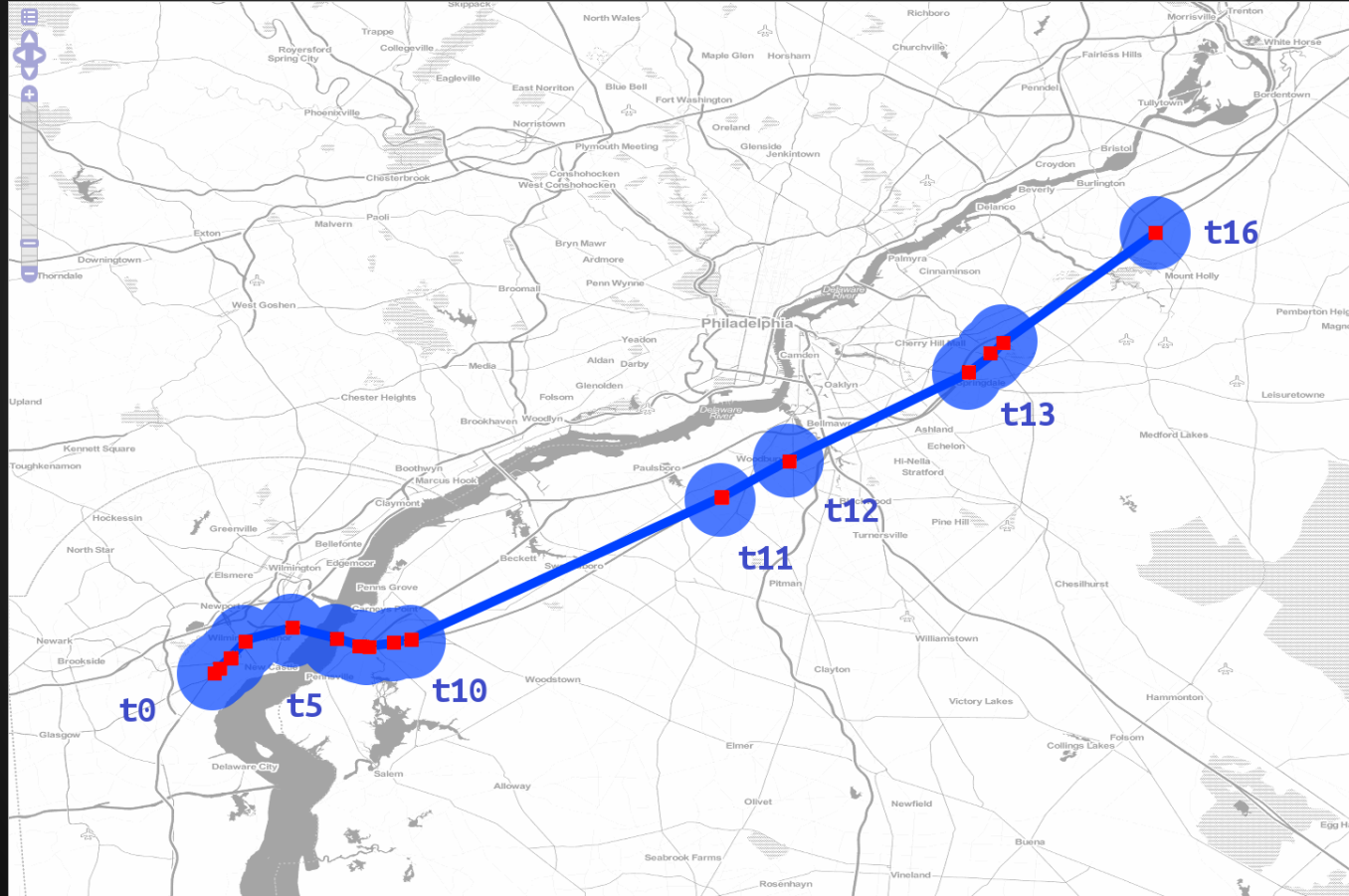
Interpolating Space-Time



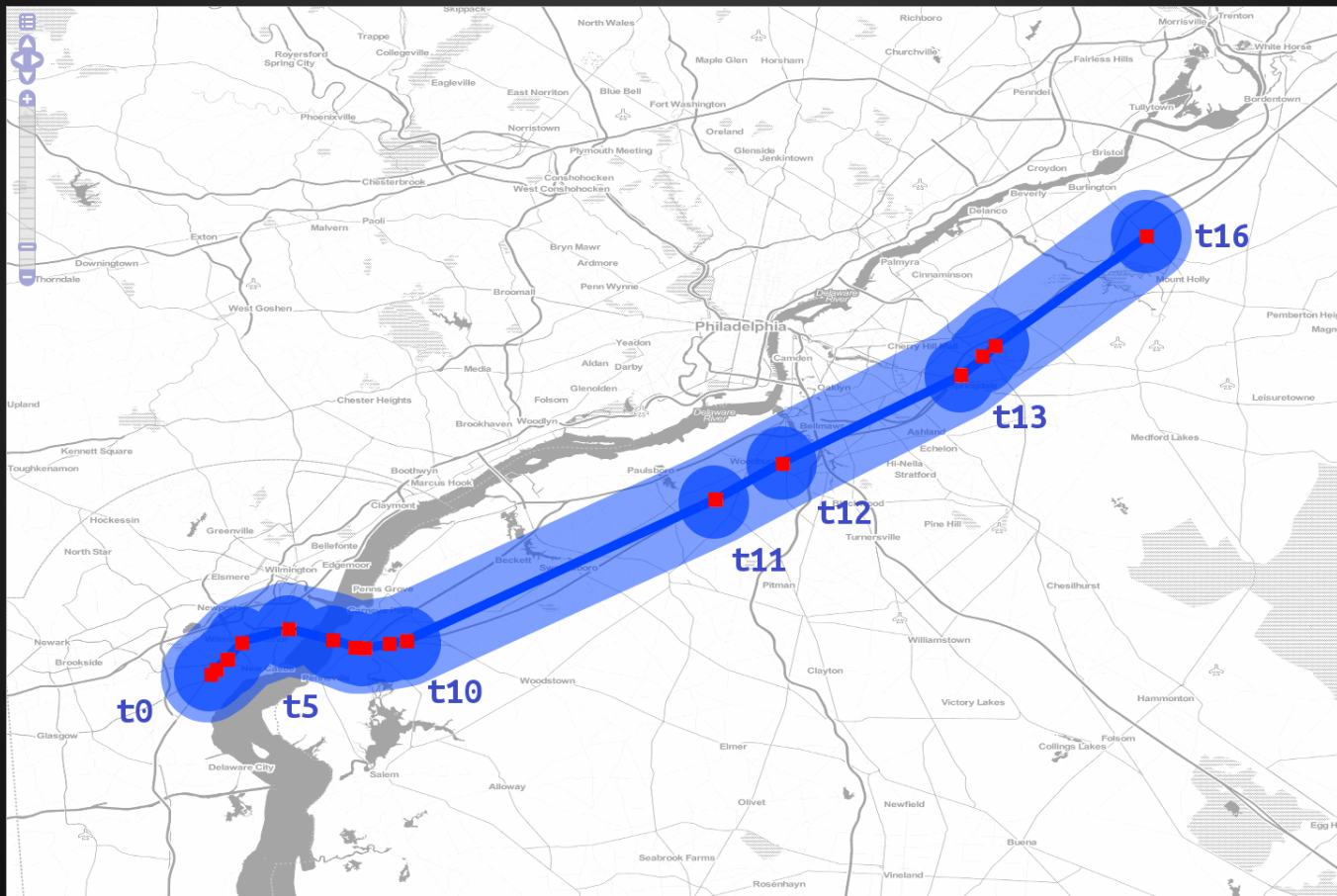
Interpolating Space-Time



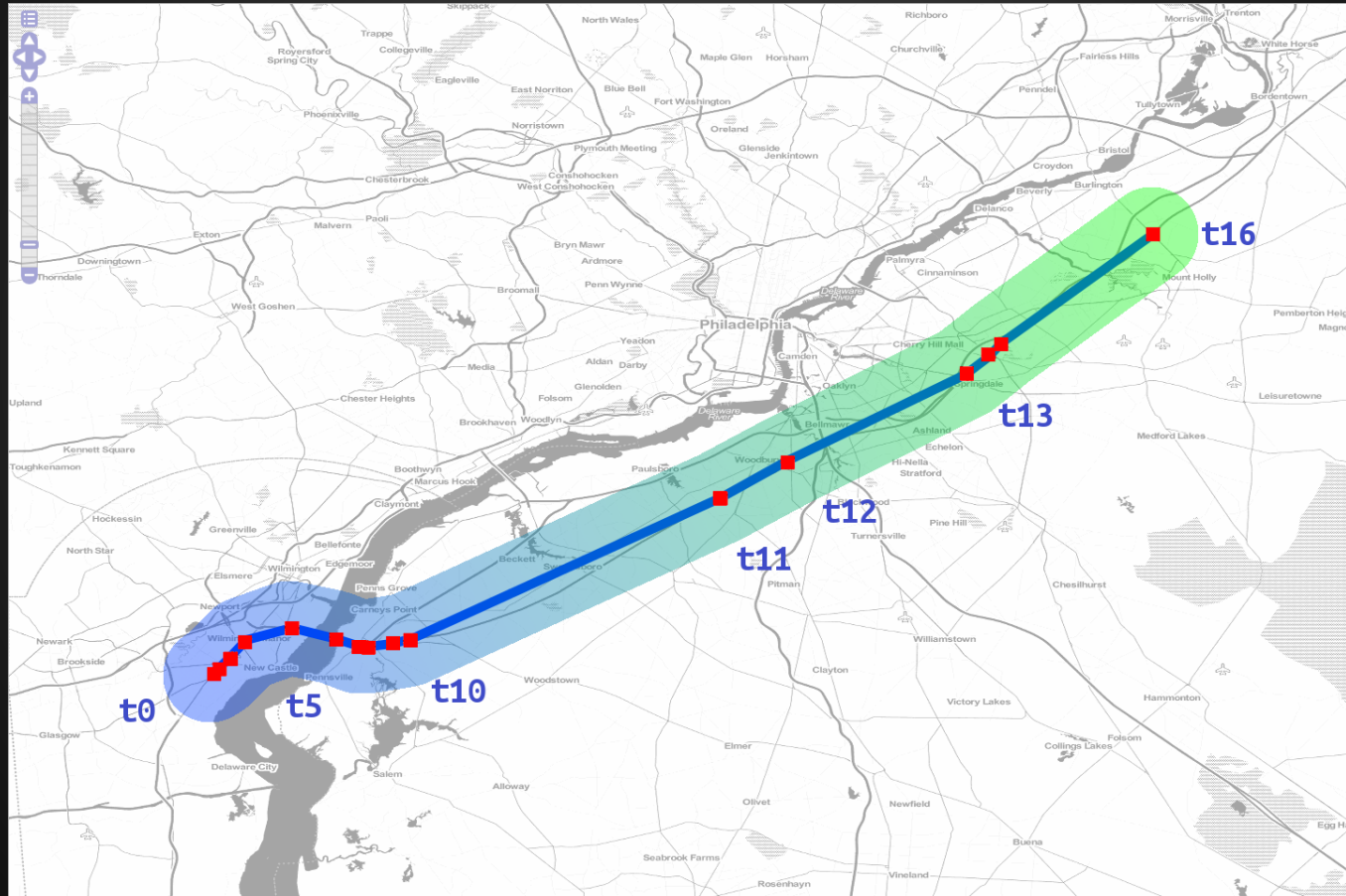
Interpolating Space-Time



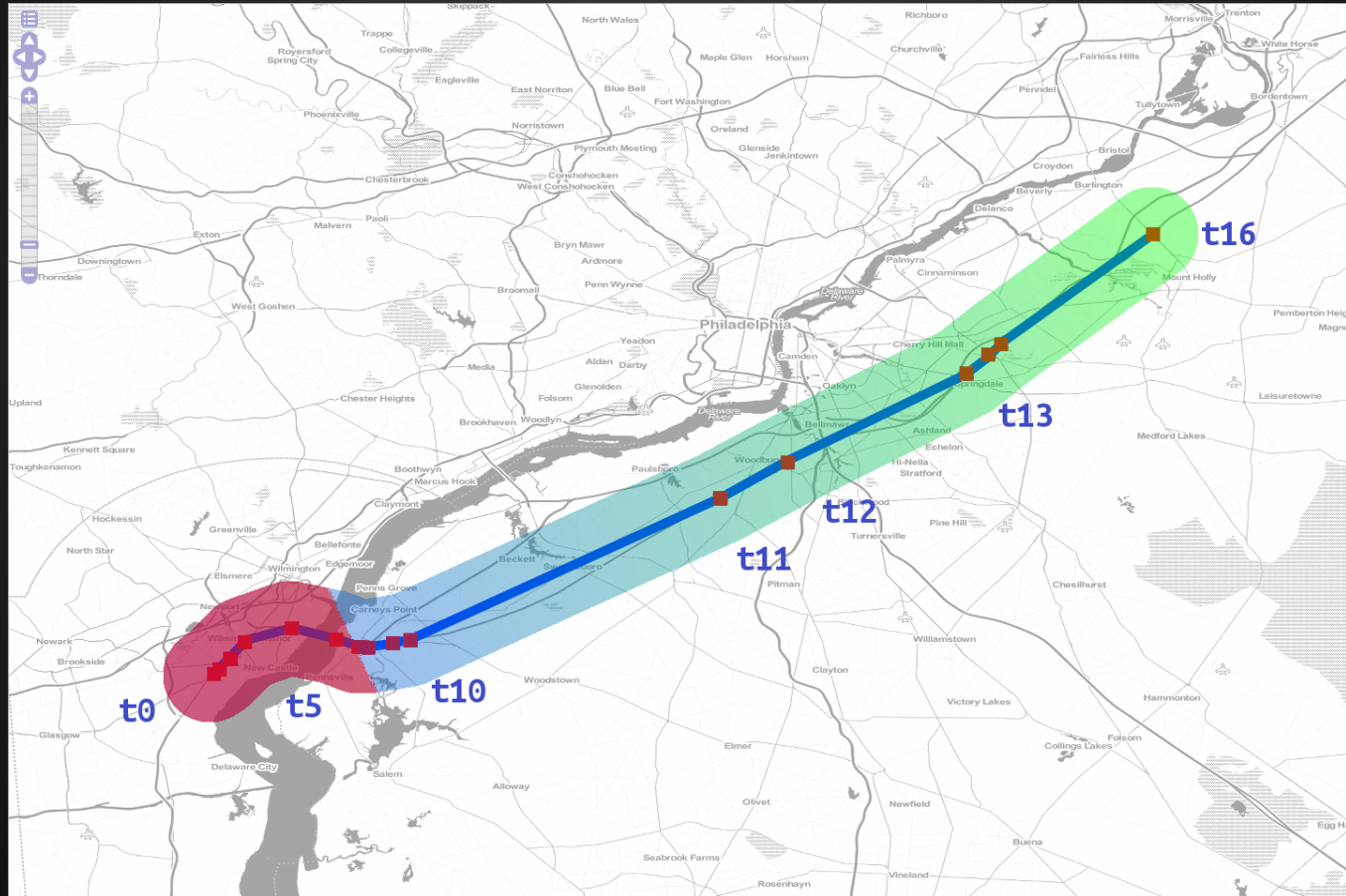
Interpolating Space-Time



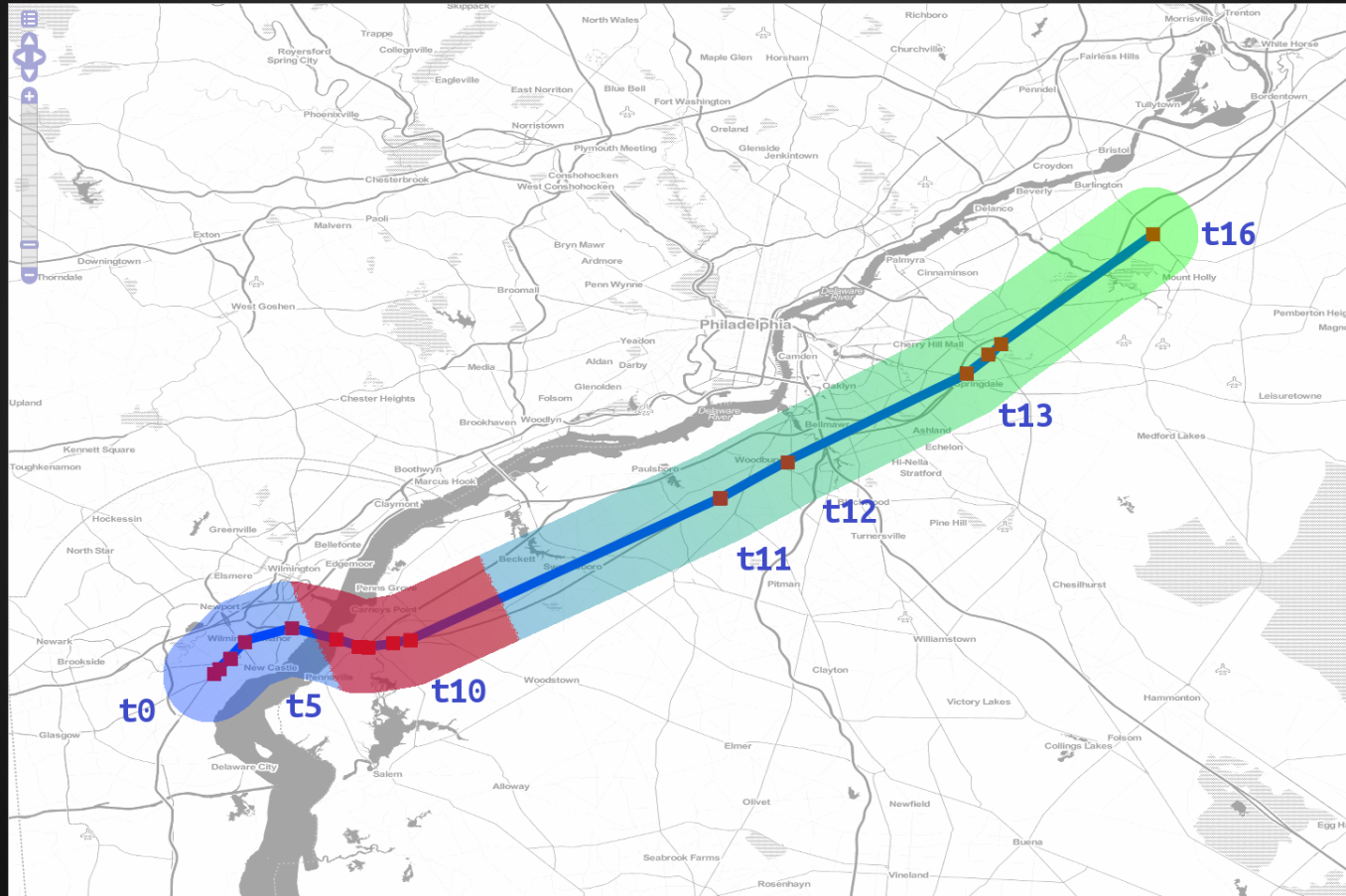
Interpolating Space-Time



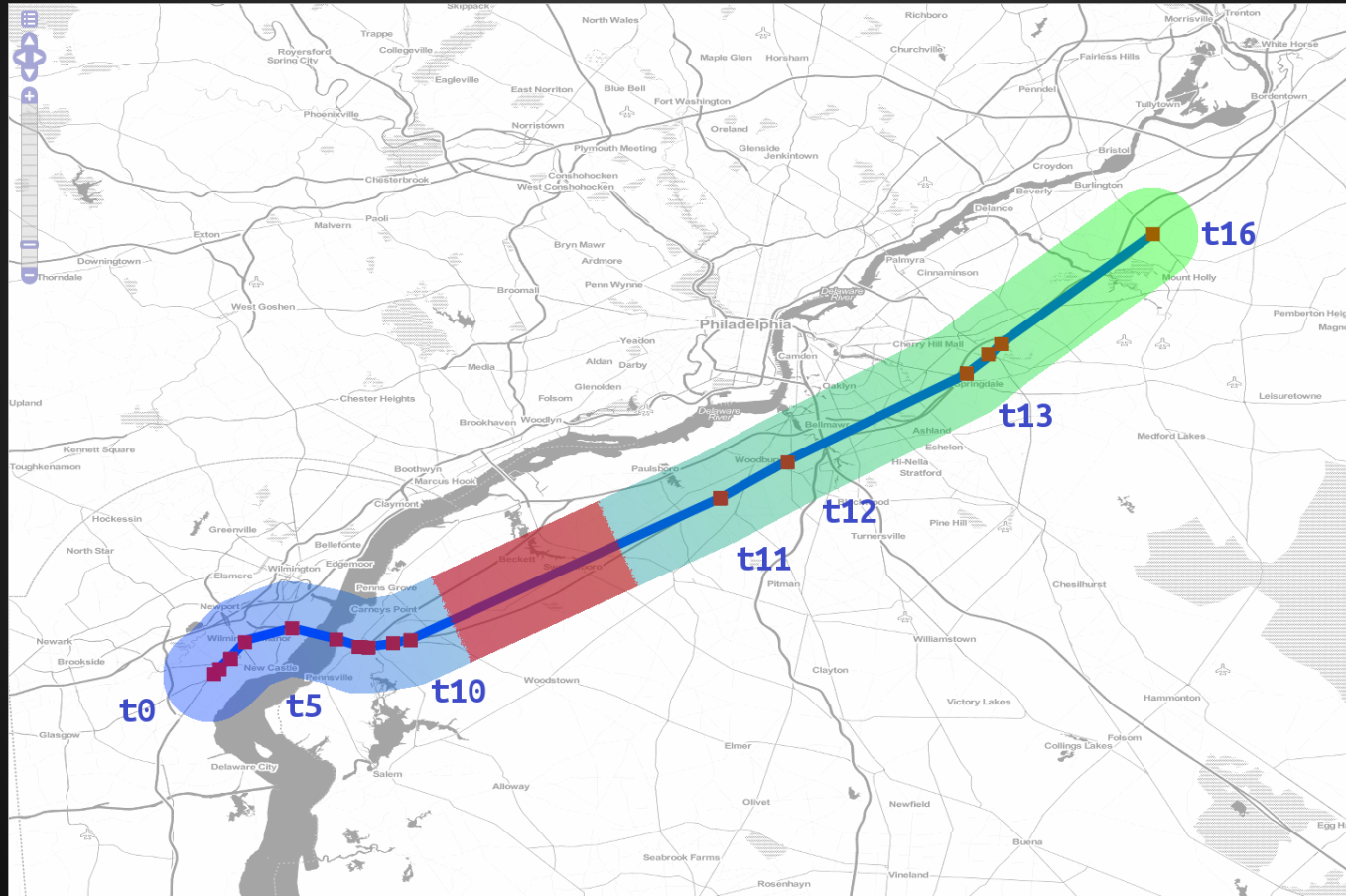
Query Planning



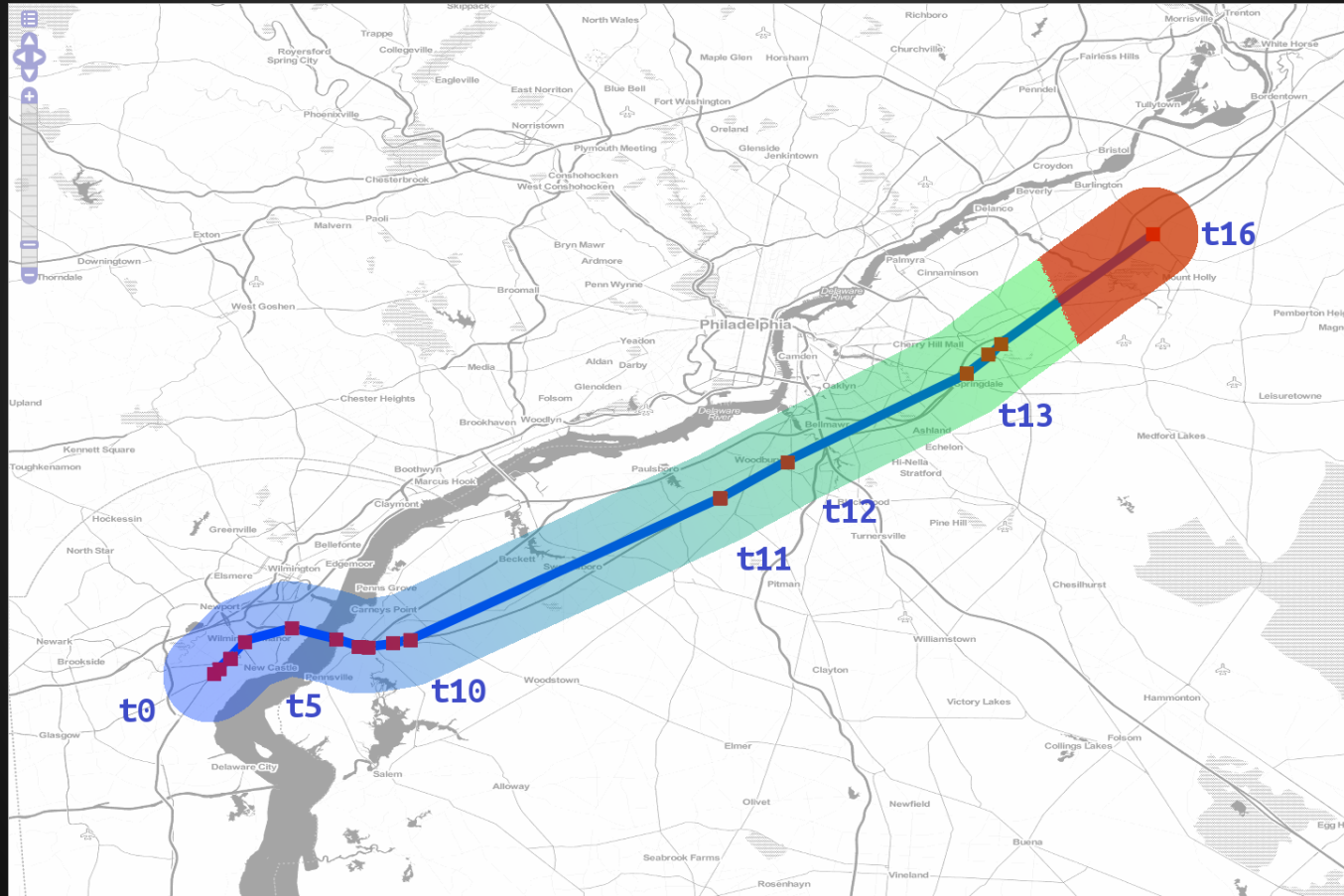
Query Planning



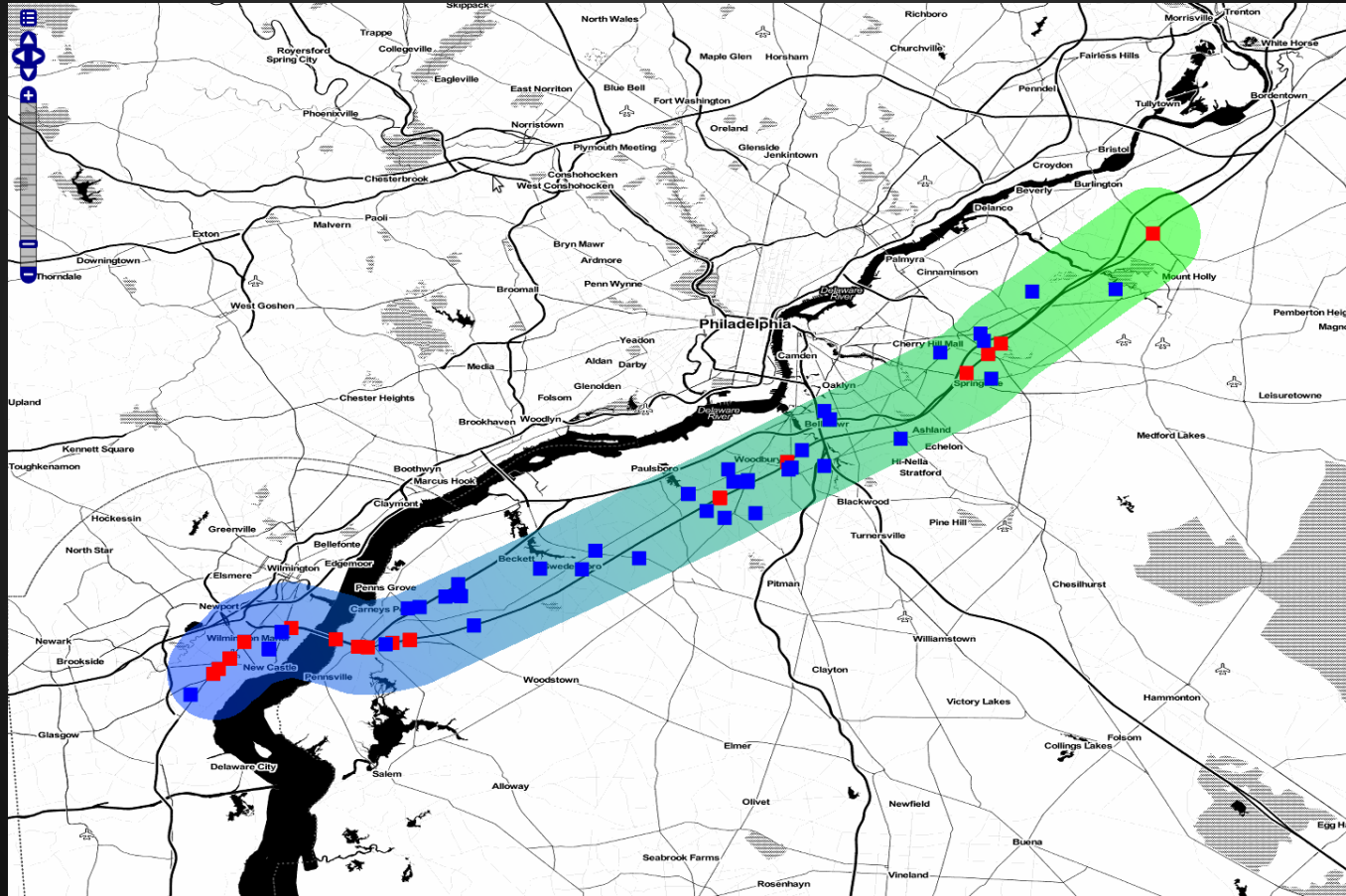
Query Planning



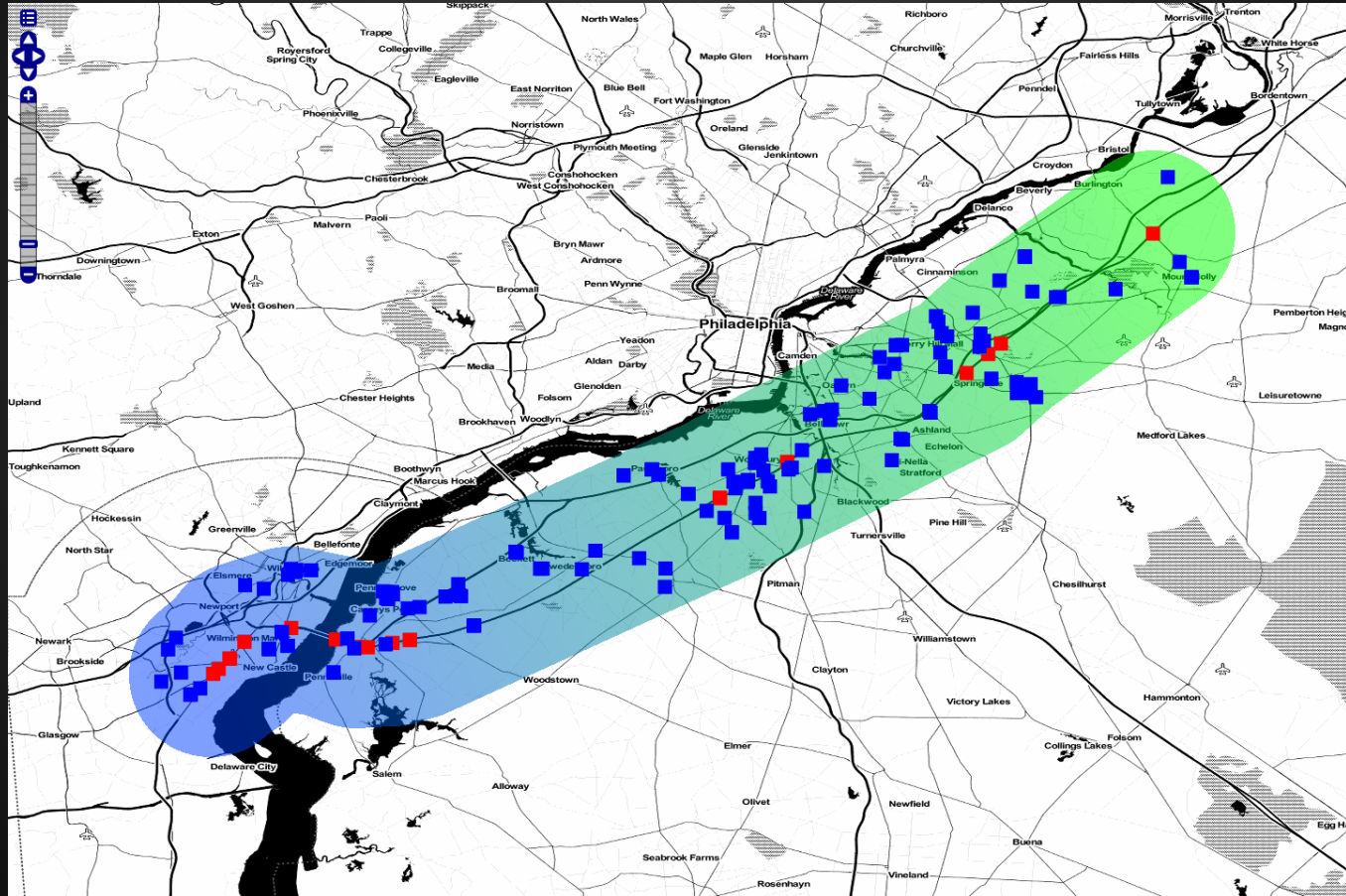
Query Planning



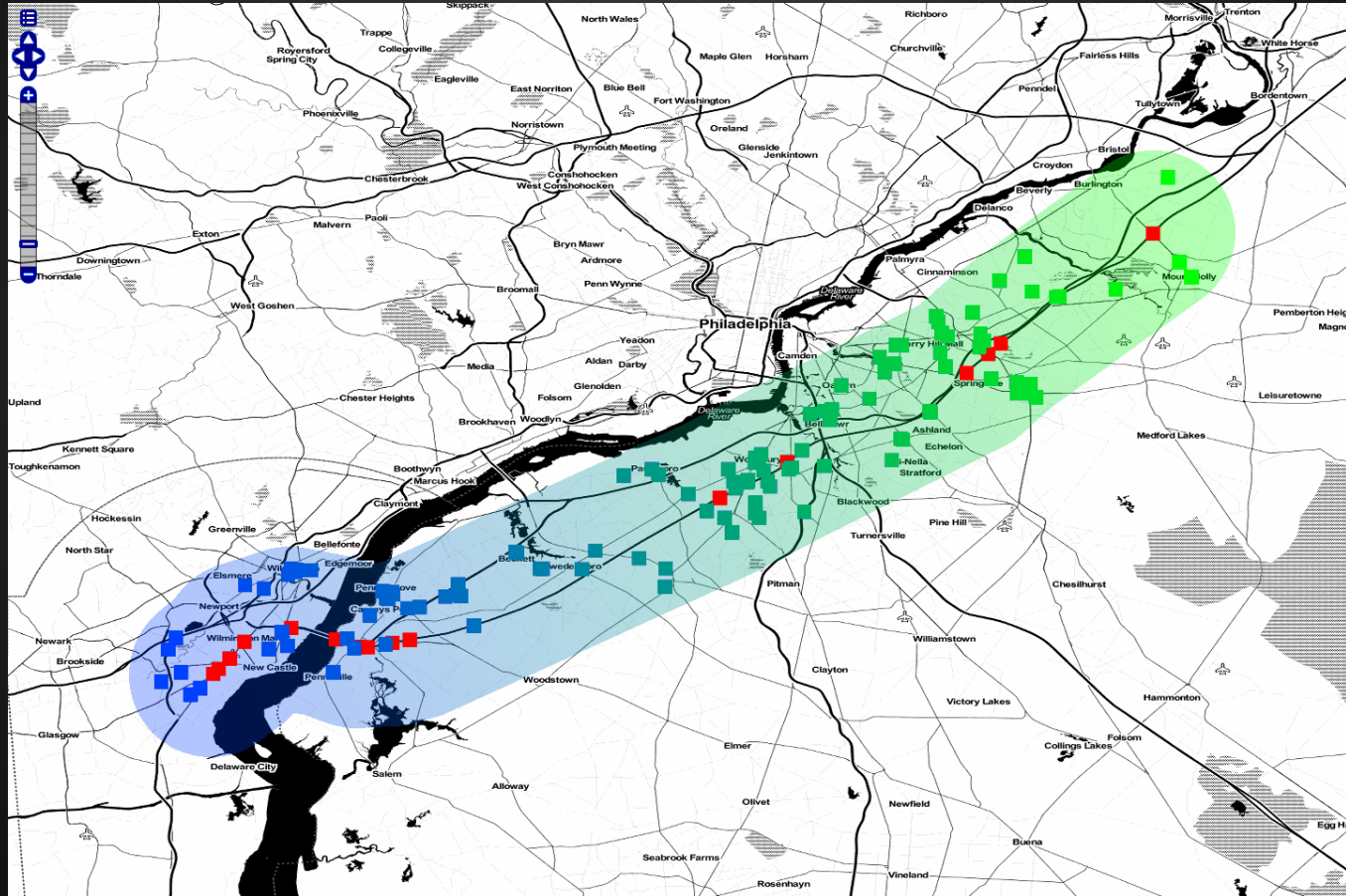
Interpolating Space-Time



Interpolating Space-Time



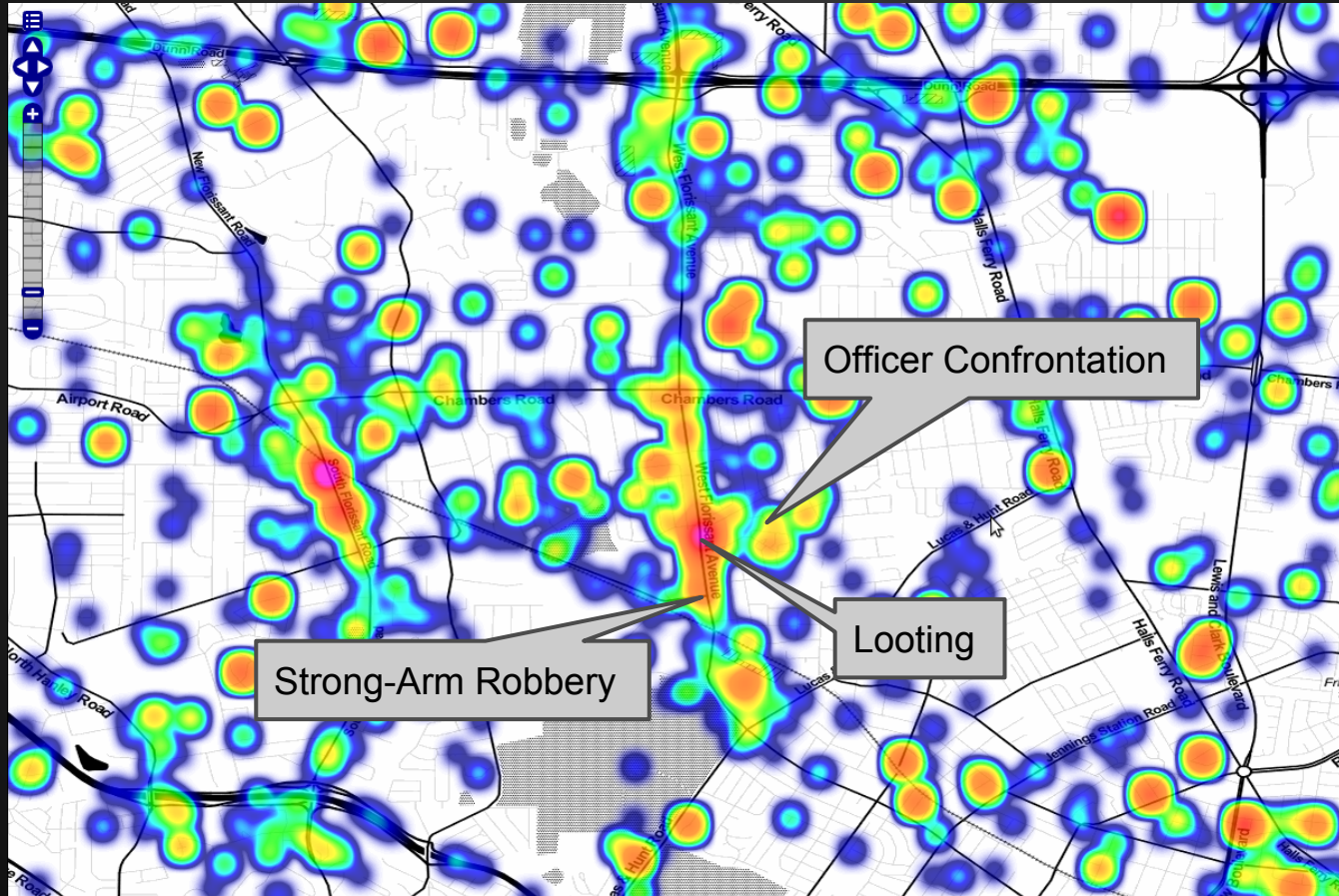
Interpolating Space-Time



Candidate Track Heuristics

- Cells Covered
- Route Cell Deviation
- $tf*idf$
- Motion Scoring
- Combined Scores

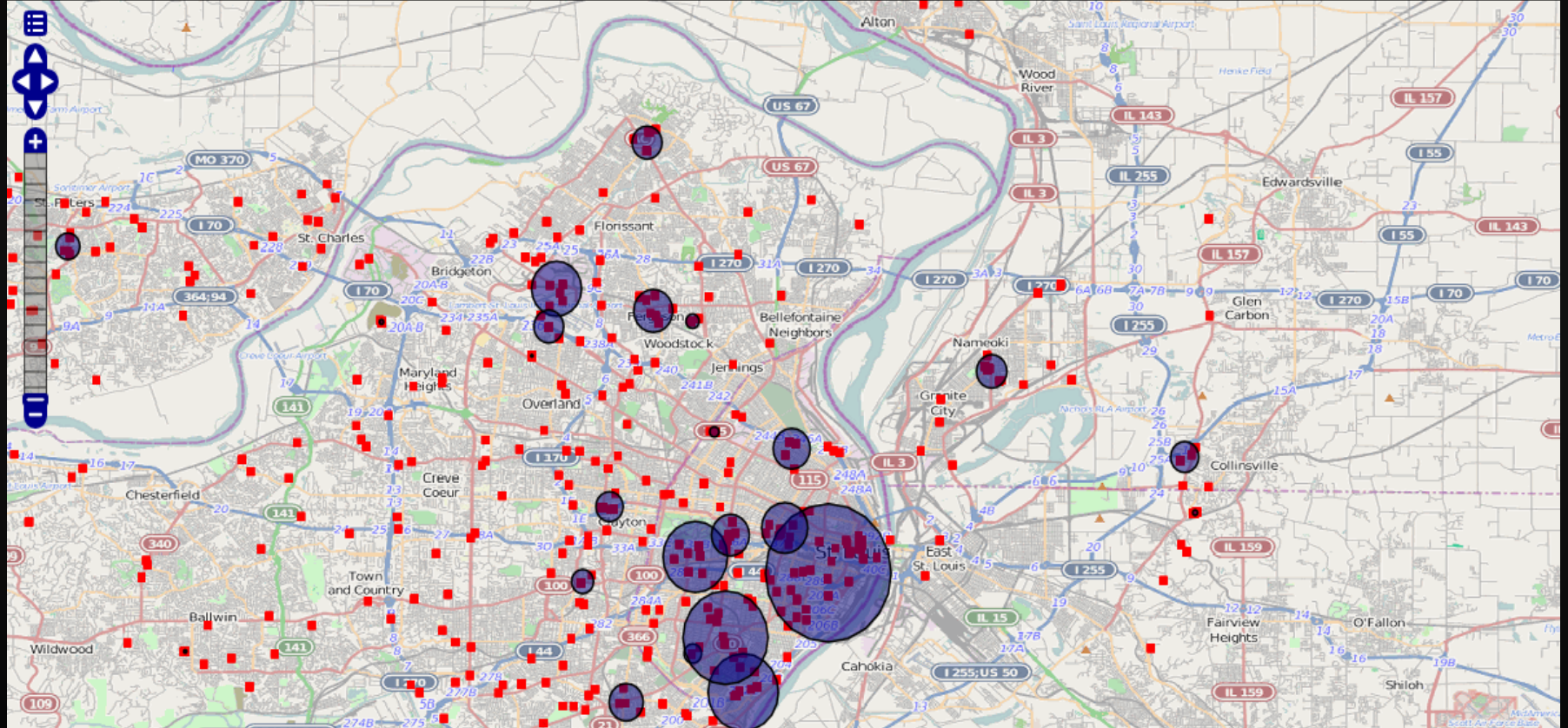
Current Events Analysis



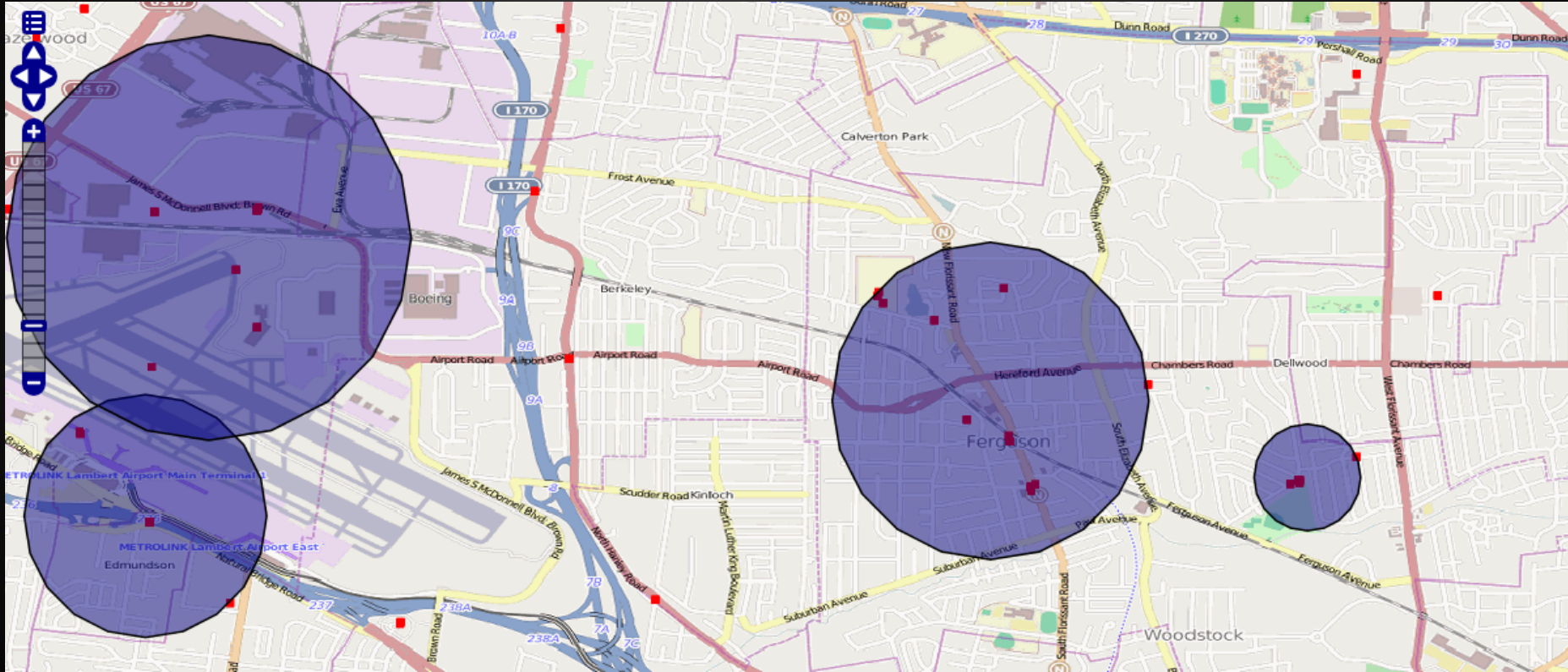
DBSCAN Clustering

- Density-based Spatial Clustering of Applications with Noise
- Clusters Algorithmically
 - No need to specify k (number of clusters)
- Arbitrarily shaped clusters

Tweet Clustering - St. Louis



Tweet Clustering - St. Louis



Apache Spark

- Real-time data analysis
- RDD (Resilient Distributed Dataset)
 - Transformations
 - Filters
 - Aggregations
- Example Time Series Output
 - Visualize via R



CQL BBOX Query

```
// Get a handle to the data store
val params = Map(
  "instanceId" -> "myinstance",
  "zookeepers" -> "zoo1,zoo2,zoo3",
  "user"       -> "username",
  "password"   -> "password",
  "tableName"  -> "geomesa_catalog")

val ds = DataStoreFinder.getDataStore(params).asInstanceOf[AccumuloDataStore]

// Construct a CQL query to filter by bounding box
val ff = CommonFactoryFinder.getFilterFactory2
val f = ff.bbox("geom", -90.32023,38.72009,-90.23957,38.77019, "EPSG:4326")
val q = new Query(feature, f)
```

GeoMesa RDD

```
val conf = new Configuration
val sconf = init(new SparkConf(true), ds)
val sc = new SparkContext(sconf)

val queryRDD = geomesa.compute.spark.GeoMesaSpark.rdd(conf, sconf, ds, query)
```

Projection & Aggregation

```
// Convert RDD[SimpleFeature] to RDD[(String, SimpleFeature)] where the first  
// element of the tuple is the date to the day resolution  
val dayAndFeature = queryRDD.mapPartitions { iter =>  
    val df = new SimpleDateFormat("yyyyMMdd")  
    val ff = CommonFactoryFinder.getFilterFactory2  
    val exp = ff.property("dtg")  
    iter.map { f => (df.format(exp.evaluate(f).asInstanceOf[java.util.Date]), f) }  
}  
  
// Aggregate and output  
val groupedByDay = dayAndFeature.groupBy { case (date, _) => date }  
val countByDay = groupedByDay.map { case (date, iter) => (date, iter.size) }  
countByDay.collect.foreach(println)
```

Time Series in R

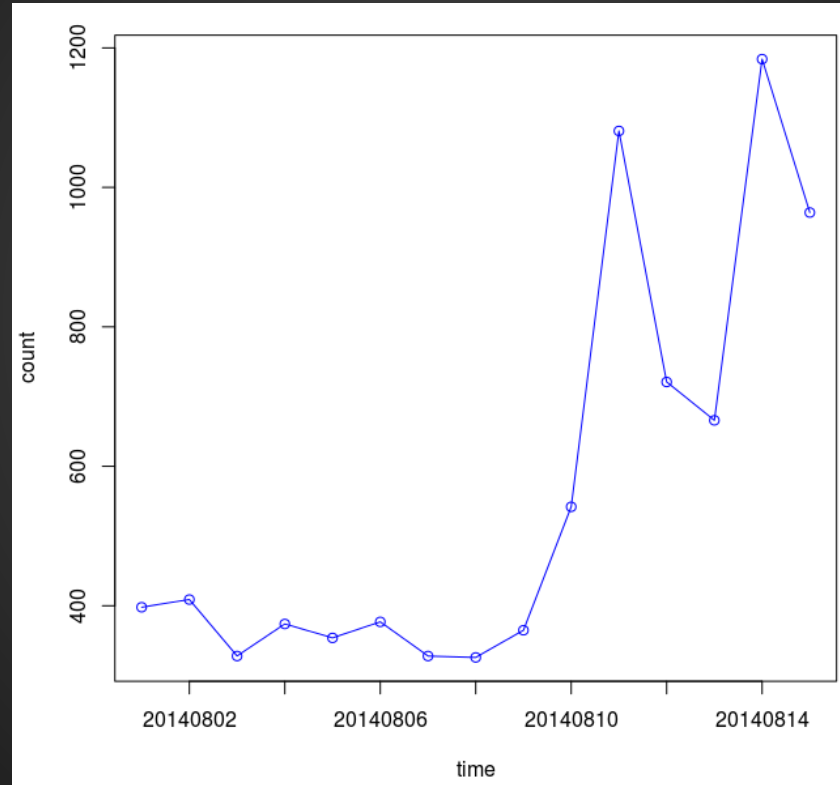
14/08/18 01:05:57 INFO SparkContext: Job finished: collect at
Runner.scala:61, took 44.154914093 s

(20140801,398)
(20140802,409)
(20140803,328)
(20140804,374)
(20140805,354)
(20140806,377)
(20140807,328)
(20140808,326)
(20140809,365)
(20140810,542)
(20140811,1081)
(20140812,721)
(20140813,666)
(20140814,1184)
(20140815,964)

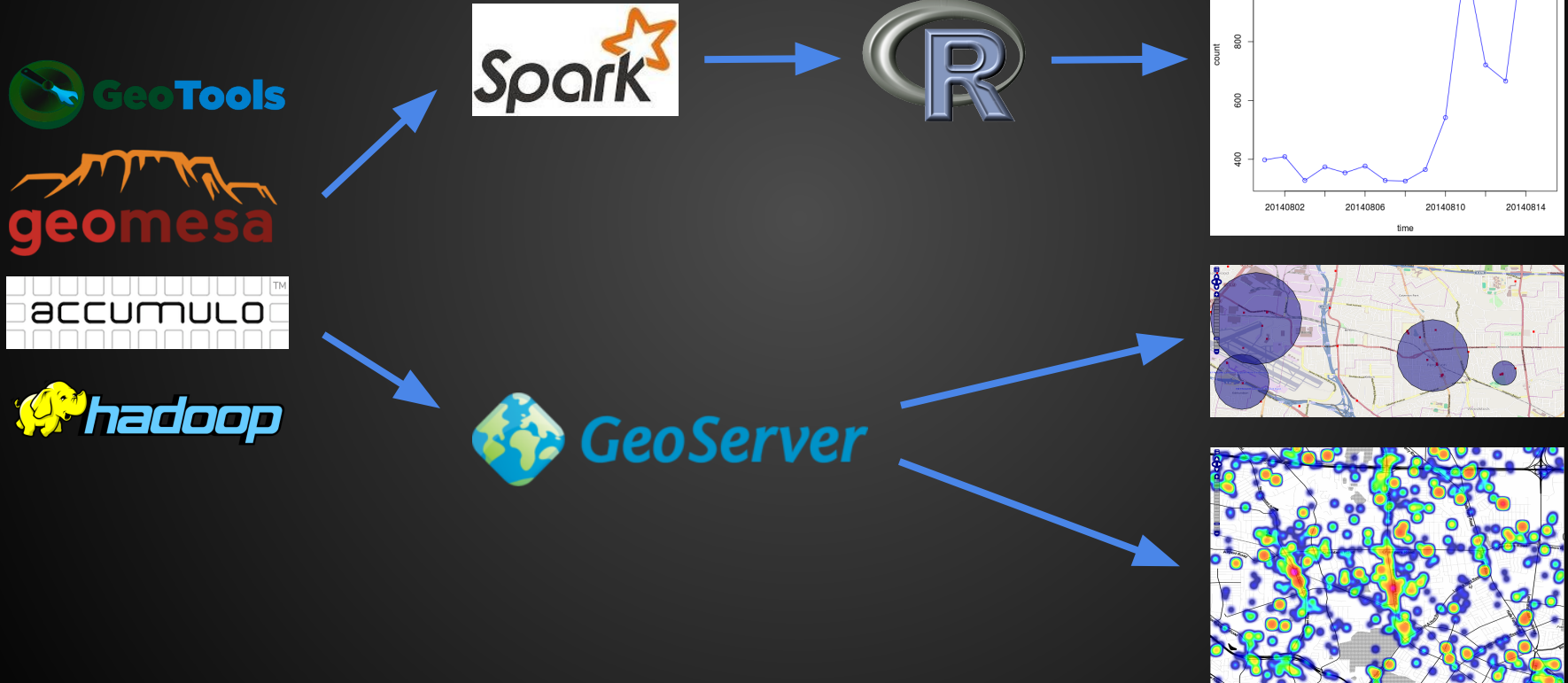
```
# read tweets
tweets = read.csv("~/Desktop/twitter.csv")

# plot tweets
plot(tweets, col="blue", type="o")
```


R Visualization



OpenSource Geospatial Analytic Pipeline



Roadmap

- Query Planning
 - Stats & Automatic Index Optimization
 - Multi-level Geospatial Indexes
- Enhance Security
- Ease Developer on-ramping
- Analytics
 - Predictive spatio-temporal analytics as a Service
 - Parallel DBSCAN
- LocationTech Community Growth

La fin du monde

Tutorials, Papers, Code, & more

<http://geomesa.org>

<http://github.com/locationtech/geomesa>

Email

geomesa-users@locationtech.org

geomesa-dev@locationtech.org

