

Rapid Analytic Development on Near Real-Time Data

Austin Heyne, CCRI
aheyne@ccri.com

The problem

Over 4,400 flights were canceled within, into or out of the United States on Wednesday, according to aviation data services company FlightAware. More than 700 flights have been canceled today, and airlines have canceled over 10,000 flights this month because of the weather.

[ABC](#); March 22, 2018

British Airways has cancelled all flights from Gatwick and Heathrow as computer problems cause disruption worldwide. [...] A spokeswoman for the airline said: "We have experienced a major IT system failure that is causing very severe disruption to our flight operations worldwide."

[The Independent](#); May 27, 2017

Ultra-large container ship CSCL JUPITER ran aground on Scheldt river bank at around 0700 LT Aug 14 at Bath, Zeeland, Netherlands, while proceeding downstream en route from Antwerp to Hamburg. [...] If hull of the giant ship is breached, dramatic situation may well turn into a nightmare.

[Maritime Bulletin](#); August 15, 2017

Then there was the NotPetya cyberattack that hit Danish shipping giant Maersk in June. The crippling attack seized the industry's attention after it cost the company \$200 million to \$300 million and led to a temporary shutdown of the largest cargo terminal in the Port of Los Angeles.

[Risk Management](#); March 1, 2018

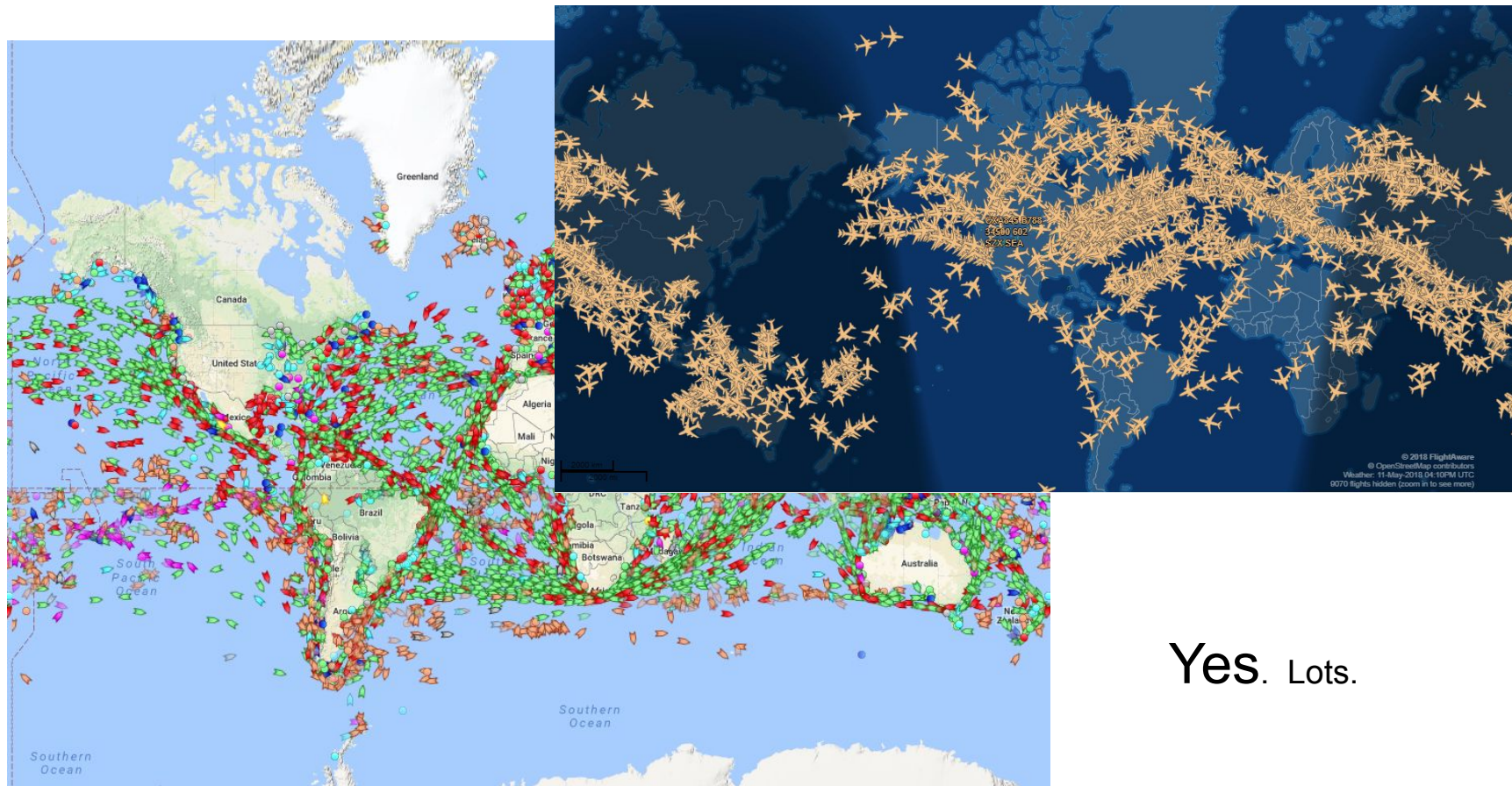
Around half of the flights in Europe on Tuesday face delays after a computer failure at the Eurocontrol center in Brussels, Belgium.

[CNBC](#); April 3, 2018

[...] a massive power outage Sunday afternoon [at Hartsfield-Jackson Atlanta International Airport] left planes and passengers stranded for hours, forced airlines to cancel more than 1,100 flights and created a logistical nightmare during the already-busy holiday travel season.

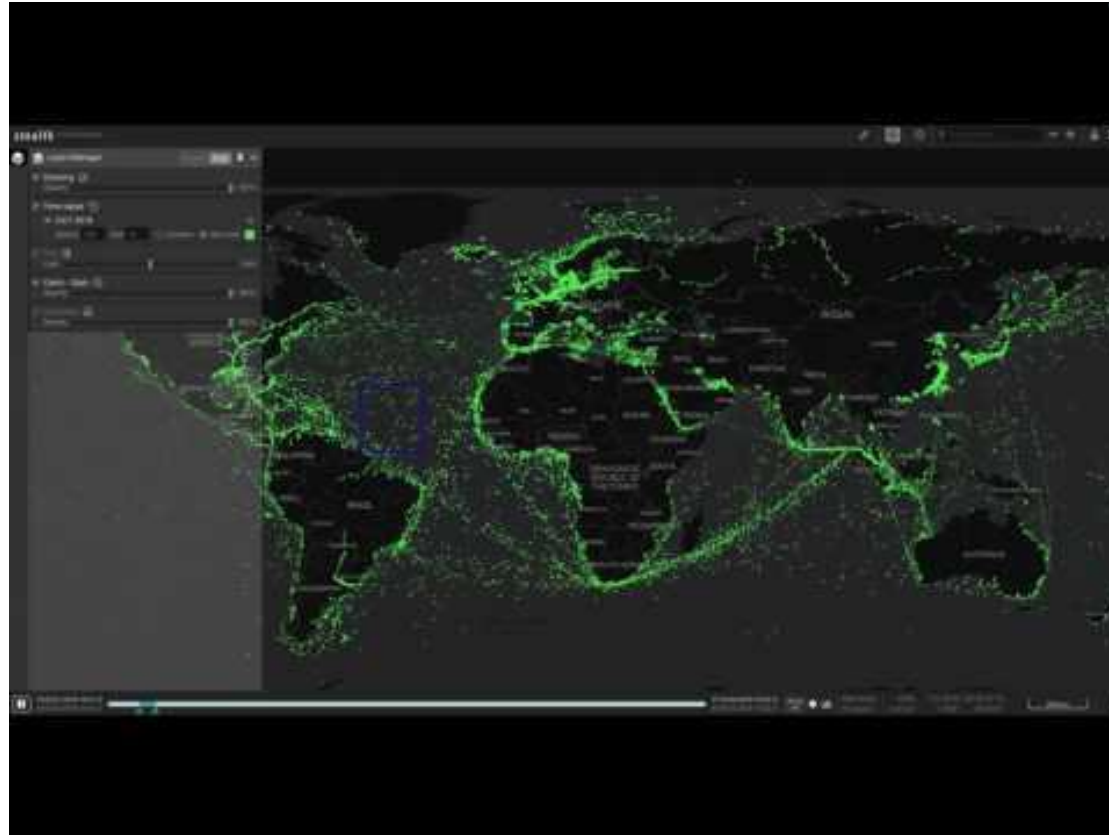
[The Atlanta Journal Constitution](#); December 18, 2017

Are there any data?



Yes. Lots.

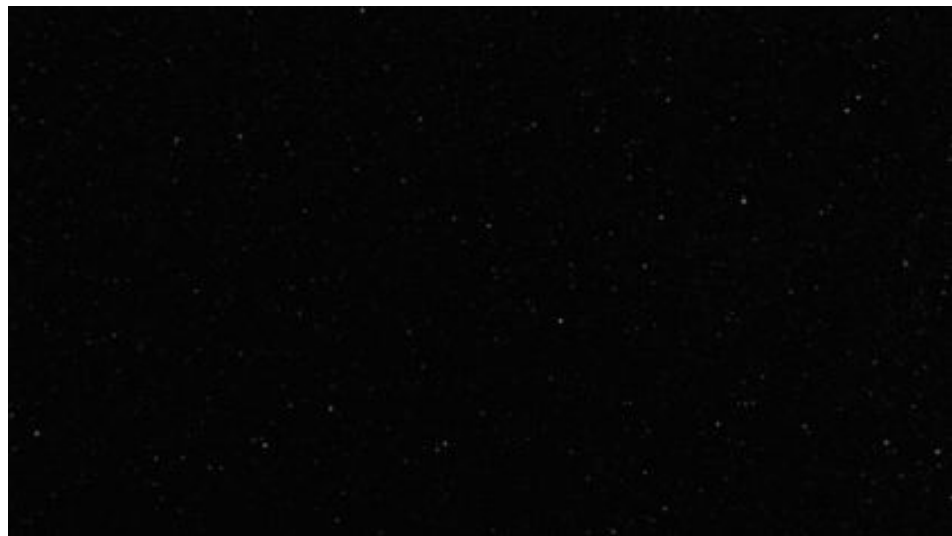
Satellite AIS



How to handle big geo-temporal data?



A suite of tools for persisting,
querying, analyzing, and
streaming spatio-temporal data at
scale... and it can SQL!



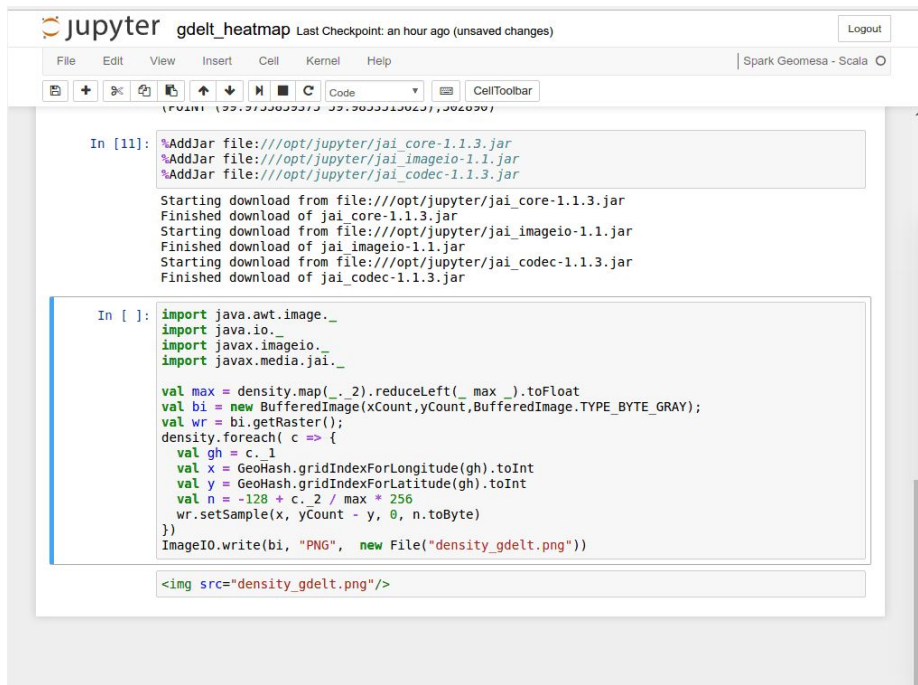
Analyst notebooks support innovation



Interactive Analysis in Notebooks

Writing (and debugging!) MapReduce / Spark jobs is slow and requires expertise.

A long development cycle for an analytic saps energy and creativity.
The answer to both is interactive 'notebook' servers like Apache Zeppelin and Jupyter (formerly iPython Notebook).



The screenshot shows the Jupyter Notebook web interface. The title bar indicates the notebook is named 'gdelt_heatmap' and shows the last checkpoint was 'an hour ago (unsaved changes)'. The interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Help) and a toolbar with icons for file operations, cell execution, and zooming. The main area displays two code cells. The first cell, labeled 'In [11]:', contains three lines of code to add JAR files: '%AddJar file:///opt/jupyter/jai_core-1.1.3.jar', '%AddJar file:///opt/jupyter/jai_imageio-1.1.jar', and '%AddJar file:///opt/jupyter/jai_codec-1.1.3.jar'. The output of this cell shows the progress of downloading each JAR file. The second cell, labeled 'In []:', contains Java code that imports necessary classes, calculates a density map, and writes it as a PNG file named 'density_gdelt.png'. The output of this cell is an HTML image tag: ''. The bottom right corner of the interface shows the user 'Spark Geomesa - Scala' and a 'Logout' button.

```
jupyter gdelt_heatmap Last Checkpoint: an hour ago (unsaved changes) Logout
File Edit View Insert Cell Kernel Help Spark Geomesa - Scala
In [11]: %AddJar file:///opt/jupyter/jai_core-1.1.3.jar
          %AddJar file:///opt/jupyter/jai_imageio-1.1.jar
          %AddJar file:///opt/jupyter/jai_codec-1.1.3.jar

Starting download from file:///opt/jupyter/jai_core-1.1.3.jar
Finished download of jai_core-1.1.3.jar
Starting download from file:///opt/jupyter/jai_imageio-1.1.jar
Finished download of jai_imageio-1.1.jar
Starting download from file:///opt/jupyter/jai_codec-1.1.3.jar
Finished download of jai_codec-1.1.3.jar

In [ ]: import java.awt.image._
import java.io._
import javax.imageio._
import javax.media.jai._

val max = density.map(_._2).reduceLeft(_ max _).toFloat
val bi = new BufferedImage(xCount,yCount,BufferedImage.TYPE_BYTE_GRAY);
val wr = bi.getRaster();
density.foreach( c => {
  val gh = c._1
  val x = GeoHash.gridIndexForLongitude(gh).toInt
  val y = GeoHash.gridIndexForLatitude(gh).toInt
  val n = -128 + c._2 / max * 256
  wr.setSample(x, yCount - y, 0, n.toByte)
})
ImageIO.write(bi, "PNG", new File("density_gdelt.png"))


```


What's missing?

We have...

- a problem
- big geo-time data
- big geo-time indexing
- interactive analyst notebook

We still need...

- a place to bring these all together
- to do it cheaply
- to do it flexibly and in a scalable way

Storage Options

- GPU RAM
- RAM
- Attached Disk (SSD, Platter)
- Cloud Blob Storage (S3, WASB)
- Cloud archive storage (Glacier)

Faster / \$\$\$



Slower / \$

Storage Options

- GPU RAM
- RAM
- Attached Disk (SSD, Platter)
- Cloud Blob Storage (S3, WASB)
- Cloud archive storage (Glacier)

Faster / \$\$\$



Slower / \$

Distributed Databases

Originally:

- Databases like Accumulo and HBase used HDFS to store data.
- HDFS needs 3-5x storage for replication and extra space.
- Compute and disk scaled together!

Distributed Databases

Originally:

- Accumulo and HBase used HDFS to store data.
- HDFS needs 3-5x storage for replication and extra space.
- Compute and disk scaled together!

Today:

- Accumulo and HBase can use Azure Blob storage and S3*
- Cloud storage scales
- Compute and disk scale separately!

* Accumulo works on Azure blob storage. Some modifications would be required to support S3. Consult your cloud professional for details, interactions, and additional suggestions.

Distributed Databases

Originally:

- Accumulo and HBase used HDFS to store data.

Today:

- Accumulo and HBase can use Azure Blob storage and S3*

Take away:

GeoMesa HBase on S3 works great!

- HDFS needs 3-5x storage for replication and extra space

- Cloud storage scales

Others have used GeoMesa Accumulo on Azure...

- Compute and disk scaled together!

- Compute and disk scale separately!

Try it today!

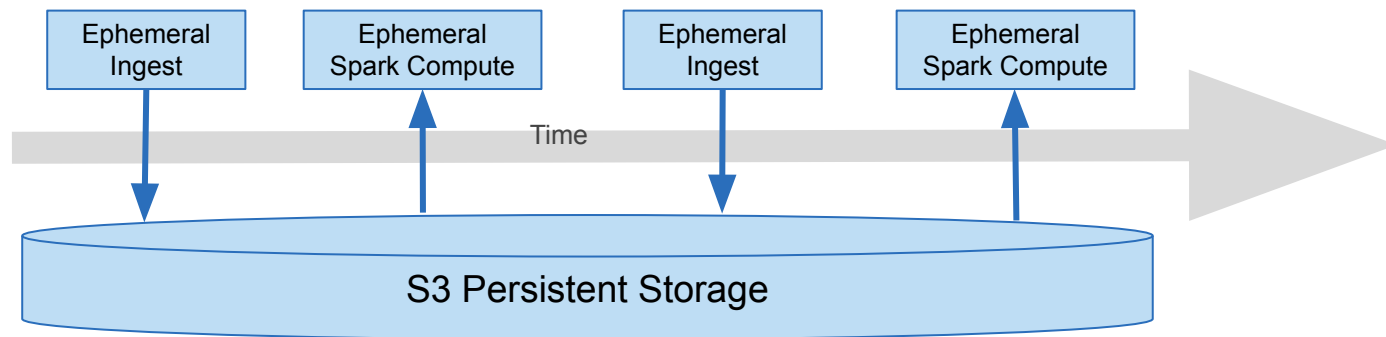
* Accumulo works on Azure blob storage. Some modifications would required to support S3. Consult your cloud professional for details, interactions, and additional suggestions.

“Ditching the Database?”

- Blobstores/Filesystems are key-value stores.
- Building a key-value store on top of a key-value store is kinda' redundant.
- GeoMesa can load and parse GDELT S3 files into Spark.
- GDELT on S3 is organized by date-keys.
- What would happen if we ‘ditched’ HBase and wrote files by space-time keys in cloud native storage?
- What parts of a ‘database’ do we really need to store/index non-updating, ascending temporal data streams?

GeoMesa FileSystem Datastore

- Serverless architecture
 - Standalone files in S3
 - Ephemeral compute for ingest and query
- Configurable partition schemes (geo + time)
- Parquet file format
 - Column-based storage (great for SQL!)
- Works great with Spark (intended for batch analysis)



Hurricane Harvey's Effect on Fuel Prices

Can oil tanker positions predict prices?

Setup

- Import GeoMesa dependency
- Create dataframe backed by GeoMesa relation
- Create SQL temporary view so we can query it

```
import org.locationtech.geomesa.spark._  
  
val dataframe = spark.read  
  .format("geomesa")  
  .option("fs.encoding", "parquet")  
  .option("fs.path", "s3a://fsds-data/")  
  .option("geomesa.feature", "AIS")  
  .load()  
dataframe.createOrReplaceTempView("fsdsAIS")
```

```
import org.locationtech.geomesa.spark._  
dataframe: org.apache.spark.sql.DataFrame = [__fid__: string, mmsi: bigint ... 36 more fields]
```

Took 1 sec. Last updated by anonymous at January 30 2018, 12:32:20 PM. (outdated)

FINISHED ▶ ⌂ ⚙

Subselect Data

- Create rough subselection of data
 - Bound by time
 - Bound by bounding box roughly around the Gulf of Mexico
- Create a new temporary view from this subselection
- Cache the data (pull into memory)

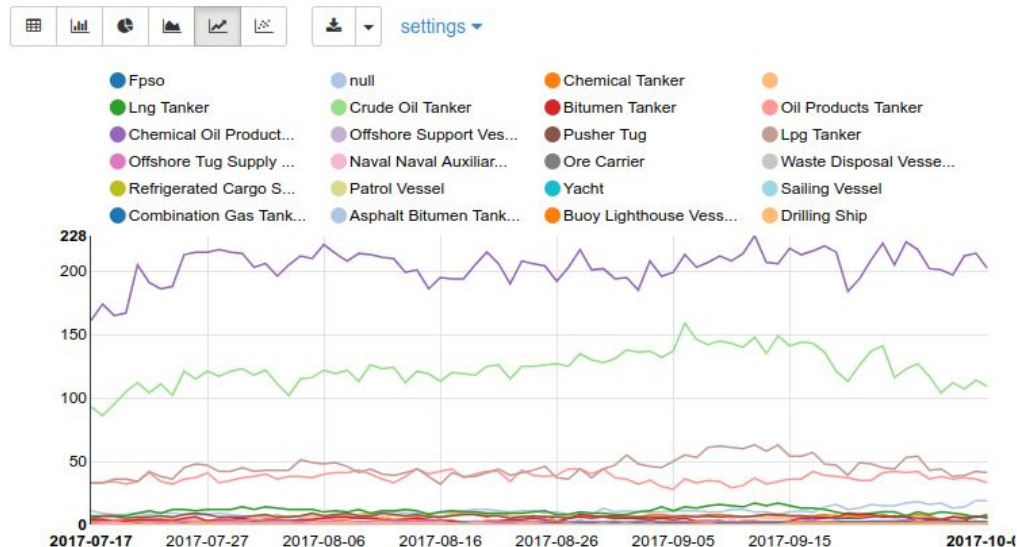
```
%sql
CREATE OR REPLACE TEMP VIEW harvey as (
  select
    *
  from
    fsdsAIS
  where
    /* Harvey formed on Aug 17 */
    /* We pad more than a month so we can later crop off some */
    /* null data and still have a clean month to examine. */
    dtg > cast('2017-07-10' as timestamp)
    and
    /* Harvey dissipated on Sept 2 */
    dtg < cast('2017-10-3' as timestamp)
    and
    /* Approximate bounding box for the Gulf of Mexico */
    st_intersects(st_makeBBOX(-100.3051, 15.5020, -73.6962, 33.8272), position)
)
```

```
%sql
cache table harvey
```

Data Exploration

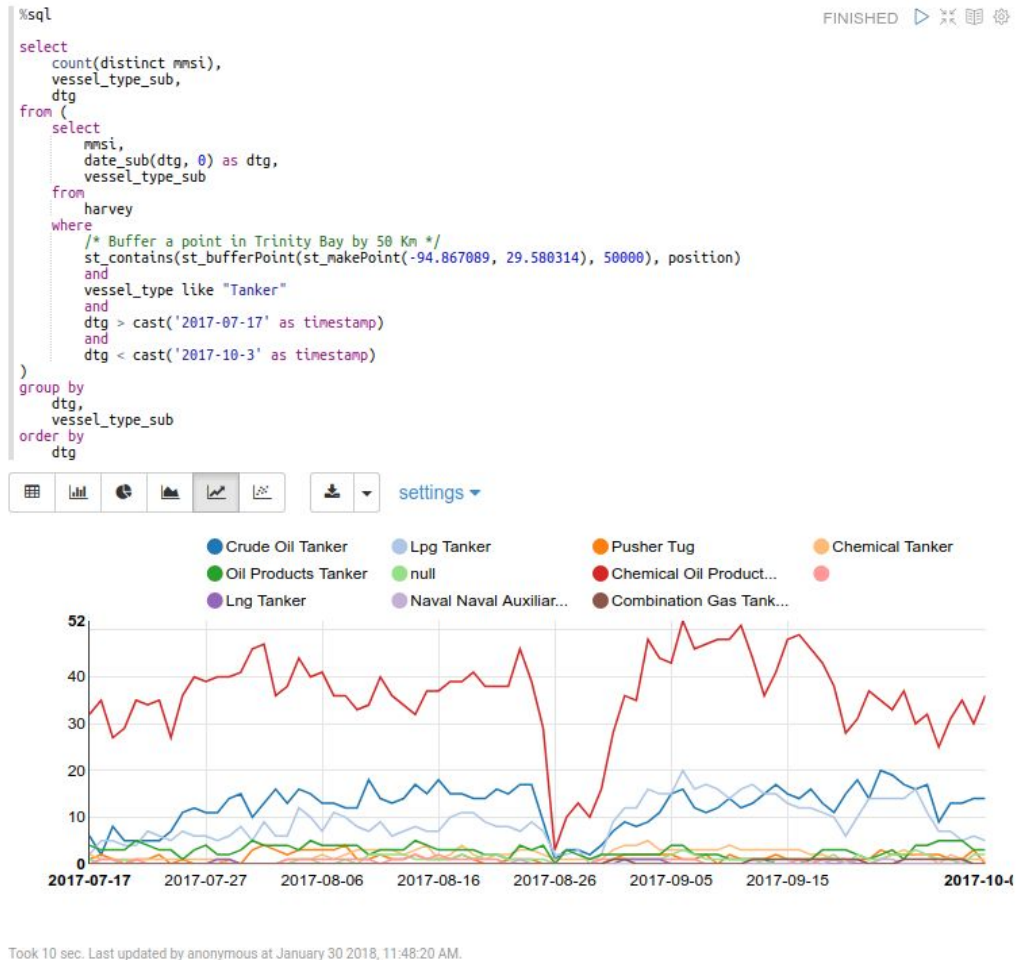
- Query for Tankers in the Gulf
- Get counts for each type of Tanker
- Group the counts by day
- Grach counts to see trends

```
%sql
select
  count(distinct mmsi),
  vessel_type_sub,
  dtg
from (
  select
    mmsi,
    date_sub(dtg, 0) as dtg,
    vessel_type_sub
  from
    harvey
  where
    vessel_type like "Tanker"
    and
    dtg > cast('2017-07-17' as timestamp)
    and
    dtg < cast('2017-10-3' as timestamp)
)
group by
  dtg,
  vessel_type_sub
order by
  dtg
```



Data Exploration

- Restrict our search to just Trinity Bay



Data Exploration

- Create a new temporary view of the number of ships in Trinity Bay

```
%sql
CREATE OR REPLACE TEMP VIEW ships as (
  /* Select the total number of tankers around houston during Harvey grouped by day */
  select
    count(distinct mmsi) as num_ships,
    /* Convert to 00:00:00 */
    date_sub(dtg, 0) as dtg_sub
  from (
    /* select ships around Houston during Harvey */
    select
      mmsi,
      dtg,
      vessel_type_sub
    from
      harvey
    where
      /* This is the Trinity Bay buffer for Houston */
      st_contains(st_bufferPoint(st_makePoint(-94.867089, 29.580314), 50000), position)
  )
  where
    /* Only get tankers */
    vessel_type like "Tanker"
  group by
    dtg_sub
)
```

FINISHED ▶ ⌂ ⚙

Extra Data

- Pull in Gas price data
 - Acquired from EIA.gov
 - Two Gas Price Indexes
 - NYH: New York Harbor
 - GC: Gulf Coast
- Create temporary view so we can analyze with SQL

```
val gasPrices = sqlContext.read
  .format("csv")
  .option("header", "true")
  .option("treatEmptyValuesAsNulls", "true")
  .option("inferSchema", "true")
  .option("mode", "DROPMALFORMED")
  .option("timestampFormat", "yyyy-MM-dd HH:mm:ss")
  .load("s3://fsds-data/csv/gasoline_spotprice_daily.csv")
gasPrices.createOrReplaceTempView("gasPrices")
```

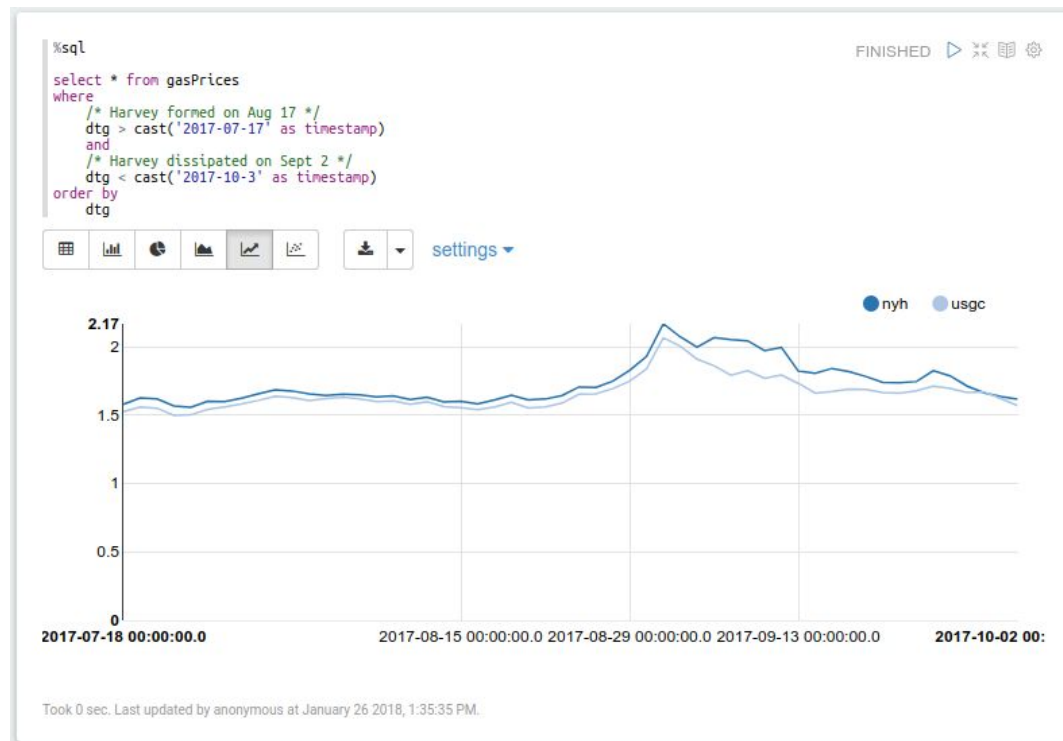
FINISHED ▶ 🔍 ⚙️

gasPrices: org.apache.spark.sql.DataFrame = [dtg: timestamp, nyh: double ... 1 more field]

Took 8 sec. Last updated by anonymous at January 30 2018, 11:45:21 AM. (outdated)

Data Exploration

- Graph data over time period of Harvey
- Notice we don't have daily values



Data Exploration

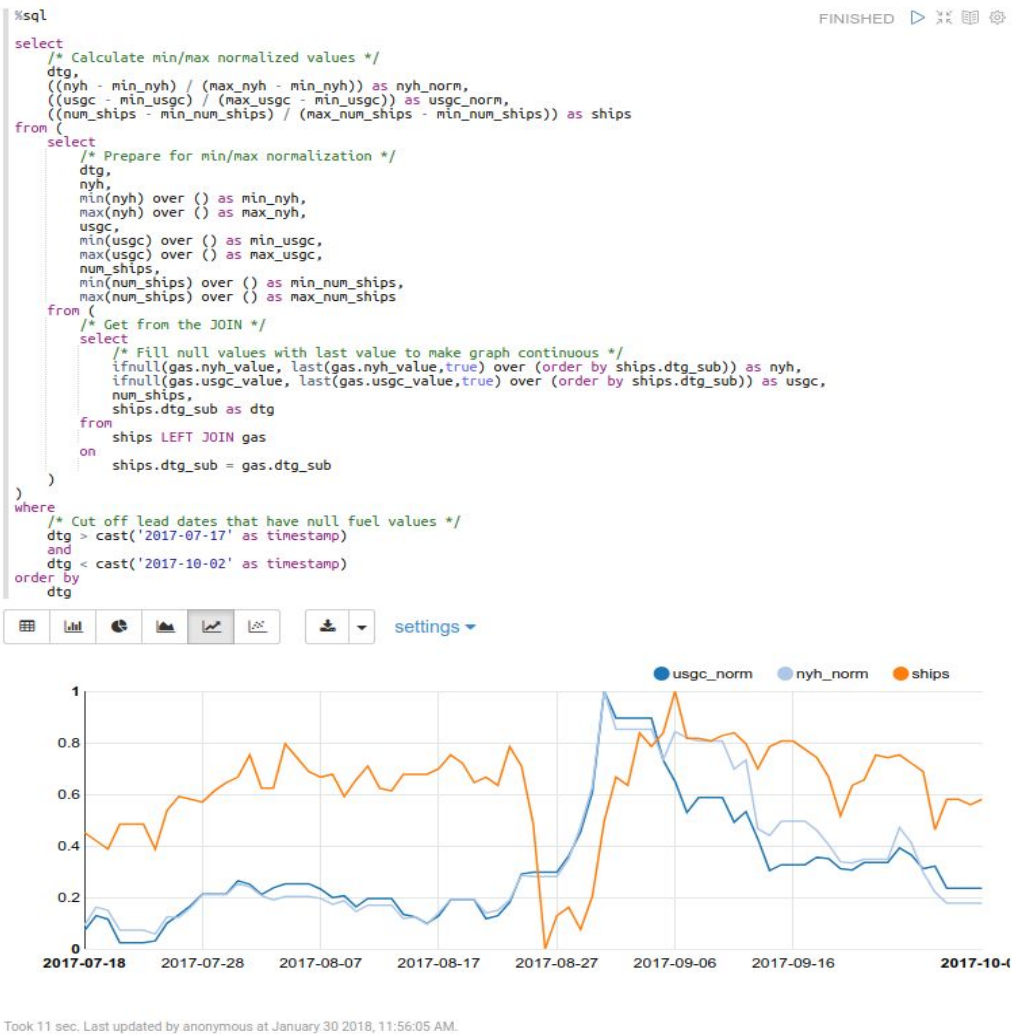
- Create temporary view of gas price data around our time of interest

```
%sql
CREATE OR REPLACE TEMP VIEW gas as (
  /* Select the Gas prices by day */
  select
    date_sub(dtg, 0) as dtg_sub,
    nyh as nyh_value,
    usgc as usgc_value
  from
    gasPrices
  where
    dtg > cast('2017-07-10' as timestamp)
    and
    dtg < cast('2017-10-02' as timestamp)
)
```

FINISHED ▶ ⌂ ⚙

Data Exploration

- Backfill the price data with the last value to give us day-continuous data
- Min/Max Normalize gas and ship counts
- Graph gas prices and ship counts together



Questions?

Find out more at <http://geomesa.org>

Connect with us on Gitter: <https://gitter.im/locationtech/geomesa>

See applications at CCRI's blog: <http://www.ccri.com/blog/>