

# SOCIAL MEDIA SEMINAR

Ryan Wesslen / Project Mosaic

Part 2: April 27, 2017

# OVERVIEW

- Recap & Json Query Example: 10 minutes
- Gnip Job Acceptance Process: 20 minutes
  - SOPHI-Gnip ingestion process
- Data Processing: 20 minutes
  - Json-to-CSV, SOPHI Twitter Datasets
- Potential Data Analysis Methods: 10 minutes

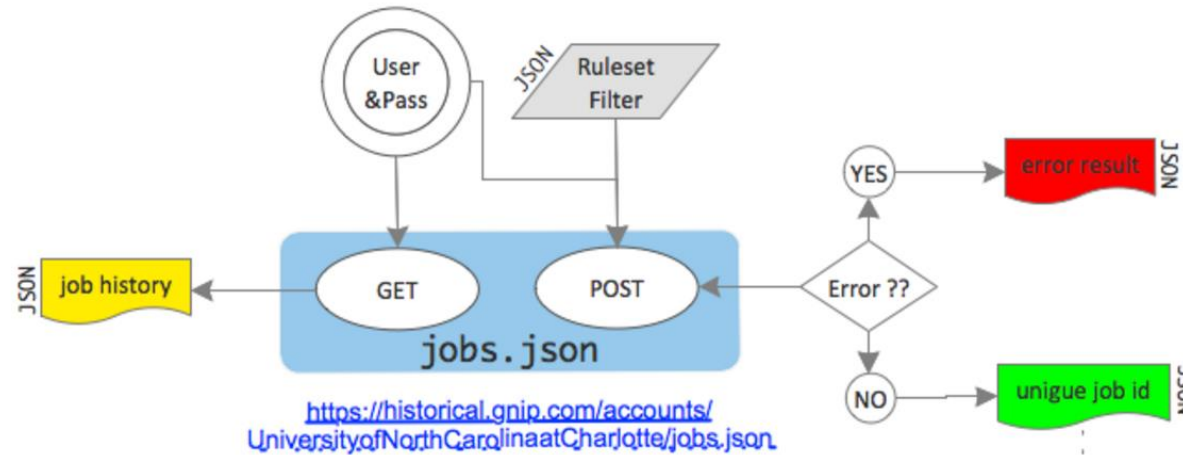
Download these slides: <http://webpages.uncc.edu/rwesslen/social-media-part2.pptx>

# JSON EXAMPLE: MARCH FOR SCIENCE

```
{
  "title": "MarchforScience",
  "publisher": "twitter",
  "streamType": "track_v2",
  "dataFormat": "activity_streams",
  "fromDate": "201704080000",
  "toDate": "201704260000",
  "serviceUsername": "twitterUsername",
  "rules": [
    {
      "value": "((March4Science) OR (March for Science) OR (Science Matters) OR #ScienceMatters OR
#ScienceMarch OR (Science March) OR (Because Science) OR #EverydayScience OR
#ScienceEmpowers OR @ScienceMarchDC OR #NoSidesInScience)",
      "tag": "marchscience"
    }
  ]
}
```

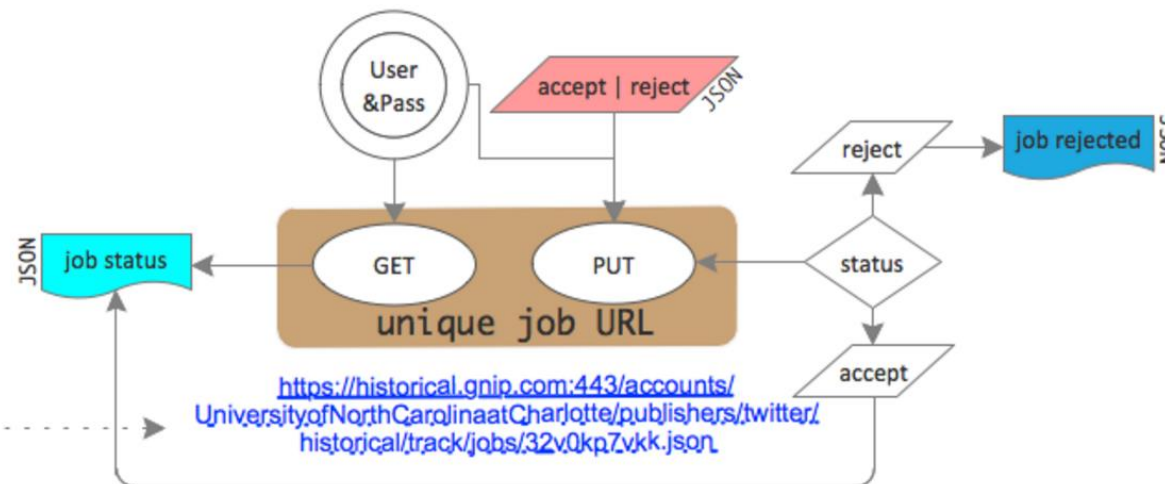
# GNIP HISTORICAL POWERTRACK PROCESS

Step 1: Submit json filter rules through command line (cURL command)



Step 2: Was the job created?

Step 3: Estimate is created



Step 4: Accept or reject the job.

# JSON RESPONSE

```
{
  "title": "MarchforScience",
  "account": "UniversityofNorthCarolinaatCharlotte",
  "publisher": "twitter",
  "streamType": "track_v2",
  "format": "activity_streams",
  "fromDate": "201704080000",
  "toDate": "201704260000",
  "requestedBy": "rwesslen@uncc.edu",
  "requestedAt": "2017-04-26T11:50:31Z",
  "status": "quoted",
  "statusMessage": "Job quoted and awaiting customer acceptance.",
  "jobURL": "https://gnip-
api.gnip.com:443/historical/powertrack/accounts/UniversityofNorthCarolinaatCharlotte/publishers/twitter/jobs/rgk4xg54c0.json"
,
  "quote": {
    "estimatedActivityCount": 1113000,
    "estimatedDurationHours": "7.0",
    "estimatedFileSizeMb": "748.39",
    "expiresAt": "2017-05-03T12:57:41Z"
  },
  "percentComplete": 0
}
```

# SOPHI INGESTION PROCESS

- Once completed, a URL (endpoint) is provided to download the files.
- The files come as **distributed** 10 min files
- Instead of users pulling these files, we have an automated process that will download the files, archive on SOPHI, and concatenate the files into one large file.

[http://support.gnip.com/apis/historical\\_api/overview.html](http://support.gnip.com/apis/historical_api/overview.html)

# JSON DATASET FILE

Example (Field Names Differ)

Tweets are stored in JSON format:

```
{ "created_at": "Wed Nov 07 04:16:18 +0000 2012",
  "id": 266031293945503744,
  "text": "Four more years. http://t.co/bAJE6Vom",
  "source": "web",
  "user": {
    "id": 813286,
    "name": "Barack Obama",
    "screen_name": "BarackObama",
    "location": "Washington, DC",
    "description": "This account is run by Organizing for Action staff.
    Tweets from the President are signed -bo.",
    "url": "http://t.co/8aJ56jcemr",
    "protected": false,
    "followers_count": 54873124,
    "friends_count": 654580,
    "listed_count": 202495,
    "created_at": "Mon Mar 05 22:08:25 +0000 2007",
    "time_zone": "Eastern Time (US & Canada)",
    "statuses_count": 10687,
    "lang": "en",
    "coordinates": null,
    "retweet_count": 756411,
    "favorite_count": 288867,
    "lang": "en"
  }
}
```

The screenshot shows a web browser window with the Hue File Browser interface. The address bar shows the URL: <https://cci-hadoop3.uncc.edu/filebrowser/view=/twitter/gnip/public/cities/Charlotte/charlotte062016.json>. The file path is displayed as Home / twitter / gnip / public / cities / Charlotte / charlotte062016.json. The file size is 50501 bytes. The preview shows a JSON array of tweet objects. The first object is for a tweet by 'Austin Caldwell' mentioning 'Kannapolis, NC'. The second object is for a tweet by 'Sandres Miller' mentioning 'UNION COUNTY, NC'.

[https://github.com/jimmoffitt/json2csv/blob/master/templates/tweet\\_standard.json](https://github.com/jimmoffitt/json2csv/blob/master/templates/tweet_standard.json)

# SOPHI DEMO

The screenshot shows the Hue File Browser interface. The browser address bar displays the URL: <https://cci-hadoopm3.uncc.edu/filebrowser/#/twitter/gnip/public>. The interface includes a navigation menu with options like Query Editors, Notebooks, Data Browsers, Workflows, Search, and Security. The main content area shows a directory listing for the path `/twitter/gnip/public`. The listing includes a search bar, action buttons (Upload, New), and a table of files and folders.

<input type="checkbox"/>	Name	Size	User	Group	Permissions	Date
<input type="checkbox"/>	↑		hdfs	hive	drwxrwxr-x	August 03, 2016 07:58 PM
<input type="checkbox"/>	·		hdfs	hdfs	drwxrwxr-x+	February 09, 2017 12:05 PM
<input type="checkbox"/>	GnipDatasets 2-9-2016.xlsx	9.1 KB	rwesslen	hdfs	-rw-r-----+	February 09, 2017 12:05 PM
<input type="checkbox"/>	business		rwesslen	hdfs	drwxr-x--+	December 19, 2016 02:10 PM
<input type="checkbox"/>	cities		rwesslen	hdfs	drwxr-x--+	December 19, 2016 03:50 PM
<input type="checkbox"/>	health		rwesslen	hdfs	drwxr-x--+	April 10, 2017 09:46 AM
<input type="checkbox"/>	politics		rwesslen	hdfs	drwxr-x--+	April 10, 2017 11:15 AM
<input type="checkbox"/>	social		rwesslen	hdfs	drwxr-x--+	March 16, 2017 03:44 PM

At the bottom of the interface, there is a pagination control showing "Show 45 of 6 items" and "Page 1 of 1".

Sophi main link: <https://sophi.uncc.edu/hue> [must be connected to UNCC network or VPN]

Twitter Datasets: <https://cci-hadoopm3.uncc.edu/filebrowser/#/twitter/gnip/public>

Twitter-Sophi Datasets Workshop: <https://webpages.uncc.edu/rwesslen/sparktwitter.pptx>



# JSON2CSV CONVERSION

## Converting Data from JSON to CSV

### json2csv

Customers often ask us about converting Tweet JSON into comma-separated values (CSV). These customers have received Twitter data from a Gnip Product such as Historical PowerTrack, 30-Day Search or Full-Archive Search, which all encode Tweets in JSON. The reasons for this are numerous, from wanting to work with a sample of the data in a spreadsheet, or needing to import the data into a relational database or legacy system. The CSV format's value is in its simplicity, and most software systems are able to import it.

To help with this process, the [json2csv app](#) is available. This tool manages the conversion of Gnip Activity Stream (AS) JSON to the comma-separated values (CSV) format.

Tweet attributes of interest are indicated by referencing a Tweet Template of choice. If the Tweet Template has an attribute it will be written to the output CSV files. If the Template does not have the attribute, it is dropped and not written. You can design your own Tweet Template, or use one of the provided [example Templates](#).

This tool works with an input folder and attempts to convert all \*.json and \*.json.gz files it finds there, writing the resulting CSV files to an output folder. This tool will work with Tweet JSON produced with both Gnip Historical PowerTrack and Search API products. Thus, this tool was designed to convert JSON Tweets in bulk.

Before deciding to perform this type of conversion, you should consider the following trade-offs:

1. JSON data are multi-dimensional, with multiple levels of nested data. However, CSVs are two dimensional. Converting from JSON to CSV means that you are sacrificing detail and flexibility in the data by either flattening it, or discarding some fields from the data.

<http://support.gnip.com/articles/json2csv.html>

Annotated instructions: [Mac](#) or [Windows](#)

Example Dataset: <https://www.dropbox.com/s/cbu1kwz2b1jlc9z/CharlotteTweets20Sample.csv?dl=0>

# HELPFUL LINKS ON GNIP DATA

- [Gnip Twitter Metadata Dictionary](#)
- [Identifying and understanding Retweets](#)
- [Filtering By Location](#)
  - [Geo Metadata](#)
- [Converting Json to CSV \(Ruby Code\)](#)
- [Consuming, Parsing and Processing Tweets with Python](#)

# DATA ANALYSIS: WORKSHOPS

- [Social Media \(Twitter, Facebook\) Data Acquisition in R](#)
- [Twitter Text Analysis with R Workshop](#)
- [R Interactive Visualizations Tutorials \(Uses Gnip Data\)](#)
- [Beer in Charlotte on Twitter Tutorial](#)
- [GitHub Repository](#)

# PARTING WORDS

- Start simple!
  - Social media is unlike traditional data sources.
  - Avoid data deluge. More data is not always better!!
- Practice, practice, practice!
  - Twitter Public API or sample datasets.
  - If you're serious, need programming skills (at least R, some Python).
- Why Twitter?
  - Avoid single platform dependency and model organism problem (Tufekci, 2012).

Project Mosaic R Faculty  
Workshop coming in June!

