

ETL and Analysis applied to Meetup

Ideas that cities and communities across
the US curious about...



Karin Brodd
Chandler McCann
Natarajan Shankar
Dan Watson

The Meetup logo, which consists of a red rounded rectangle with a white horizontal band in the center. The word "meetup" is written in a black, lowercase, sans-serif font within the white band. A small registered trademark symbol (®) is located to the upper right of the logo.

meetup

Project Scope

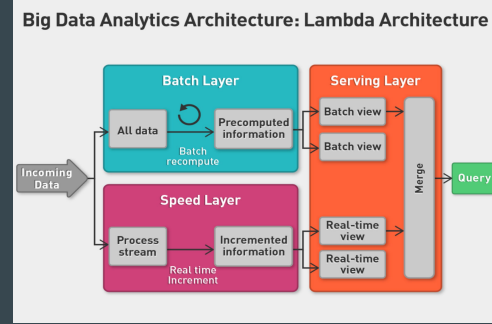
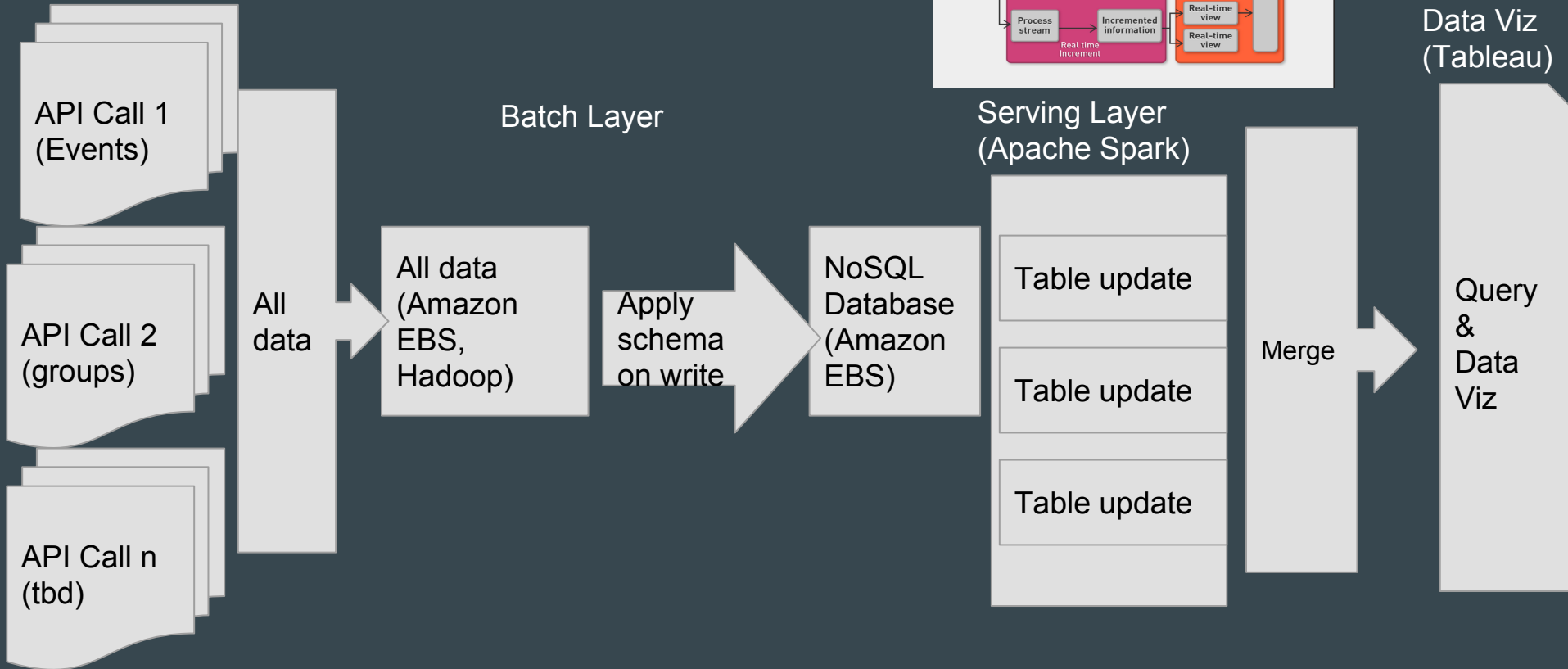
- Apply the ETL concepts learned in this class
 - Pick a provider of a live data stream
 - Establish an architecture for the end to data flow and implement
 - Refine the chosen architecture through issues encountered
 - With a working implementation:
 - Extract information about the curiosities and interests in geographical regions and subgroups (of practical scope)
 - Which communities are relatively more active throughout the US?
 - Are similar groups more active in one geographic region than in another?
 - What possibly accounts for these differences?
- Project the architecture and implementation towards a future form!



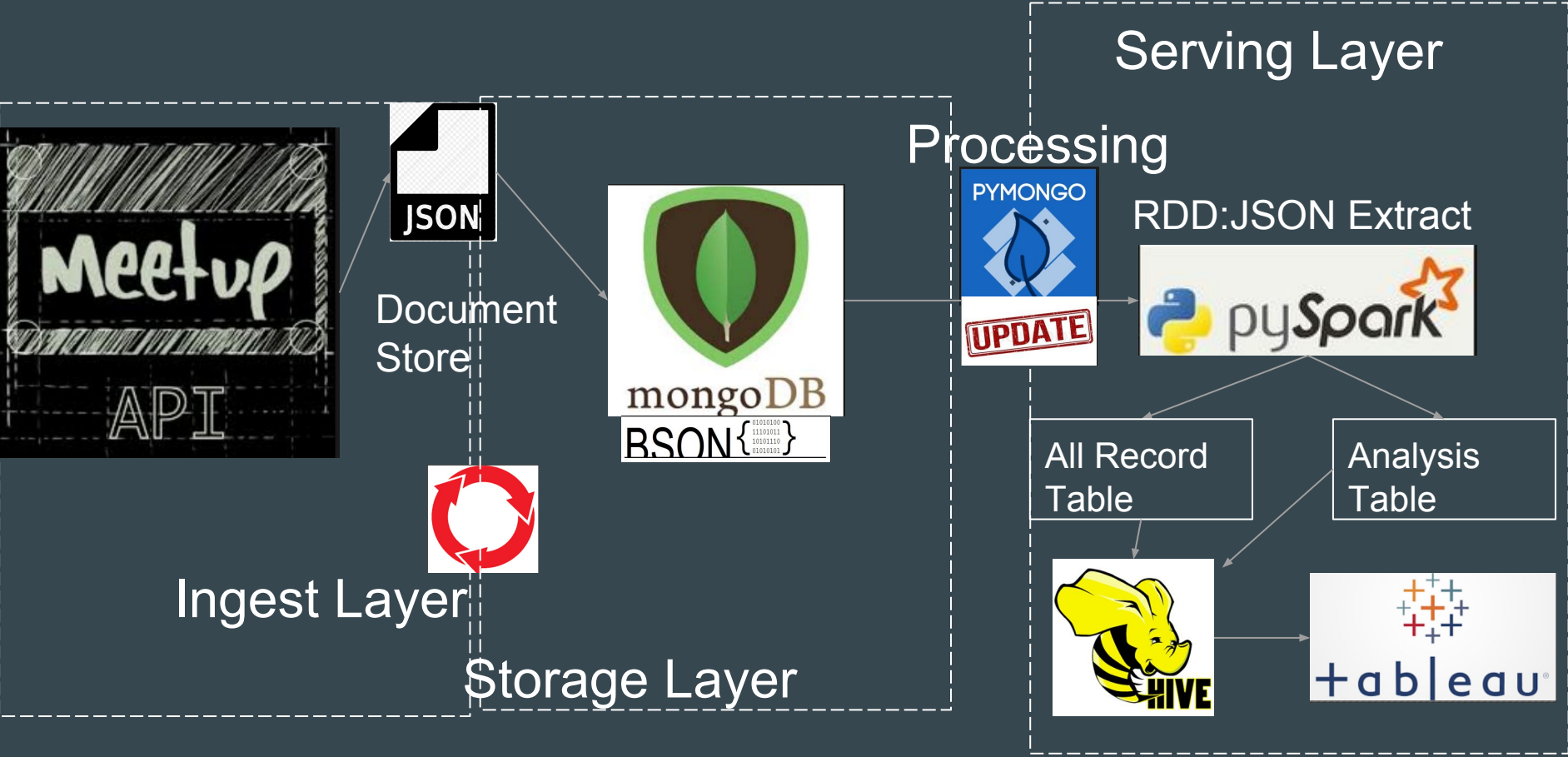
Ingestion of live data stream - Data Source : Meetup.com

- Why meetup.com?
 - Real data source, streamed
 - Meetup topics reflect public inclination and sentiment
 - Measure of frequency of similar meetings across the nation will augment and outperform national polls
 - Meetup.com offers diverse API
 - Data diversity via events and get_events API
 - Meets velocity requirements, many hundred JSON records per query
- Data Ingestion Layer
 - Preprocessing
 - Segment the incoming stream into individual JSON records
 - Built-in filtering: Write only the clean records to Database
 - Dirty records are dropped prior to database write

Initial proposed architecture... based on Classic Lambda



Refined Architecture and Data Flow



Tradeoffs, Design Choices, Rationale

- **Limit the number of spout types**
 - Design for multiple types of data spouts but focus on meetup.com get_event()
- **Use in stream schema application and data cleanup**
 - Chosen data spout provides well formatted data
 - Schema on Write is not necessary, JSON is already well formed
 - Drop bad JSON in stream
- **Choose MongoDB over S3**
 - Need a self managed DIY app rather than a managed service
 - Couple the DB very closely with the API process
 - Document store in Mongo well suited for JSON
- **Write persistent Hive tables from Pyspark**
 - Data manipulation in Pyspark with data frames, has long term machine learning abilities
 - Native hive context in Spark allows easy save of Hive tables
 - Easy access to Hive from Tableau
 - Trouble getting the spark SQL thrift server to work

The long road from Mongo to Spark :

Advanced Analytics, Complex Computation

Versatile



*Cumbersome
lynch pin*



PYMONGO_SPARK
on top of
Mongo-Hadoop

Convenient



```
import pymongo_spark  
pymongo_spark.activate()
```



- MongoDB Conn
- mongoRDD(conn)
- JSON Extract
- DF manipulation
- Scale

Easily connect to a mongoDB with a mongoRDD

```
conn = "mongodb://meetup_user:%40B%24%40rp1Dy%25C3@ec2-54-144-14-181.c
rdd = sc.mongoRDD(conn)
2/04 02:43:47 WARN SampleSplitter: Not enough documents for more than
new_rdd = rdd.map(lambda x: dict([(i, x[i]) for i in x if i != '_id']))
                ).map(lambda x: json.dumps(x, ensure_ascii=False).enco
                ).map(lambda x: "".join(x.split("\n")))
)
df = sqlContext.jsonRDD(new_rdd)
```

Schema automatically inferred from JSON

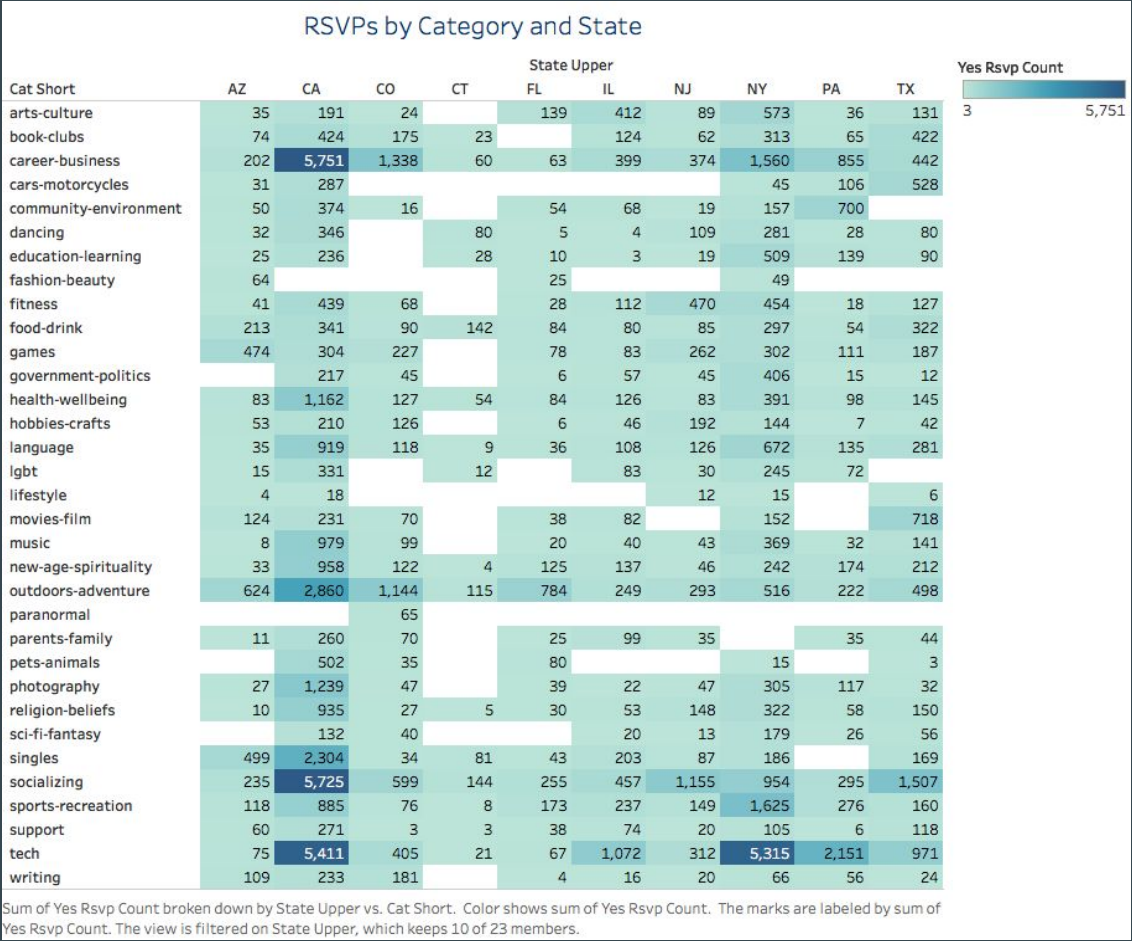
```
[>>> df.printSchema()
root
 |-- created: long (nullable = true)
 |-- description: string (nullable = true)
 |-- distance: double (nullable = true)
 |-- duration: long (nullable = true)
 |-- event_url: string (nullable = true)
 |-- fee: struct (nullable = true)
 |   |-- accepts: string (nullable = true)
 |   |-- amount: double (nullable = true)
 |   |-- currency: string (nullable = true)
 |   |-- description: string (nullable = true)
 |   |-- label: string (nullable = true)
 |   |-- required: string (nullable = true)
 |-- group: struct (nullable = true)
```

Create columns from nested
JSON fields easily



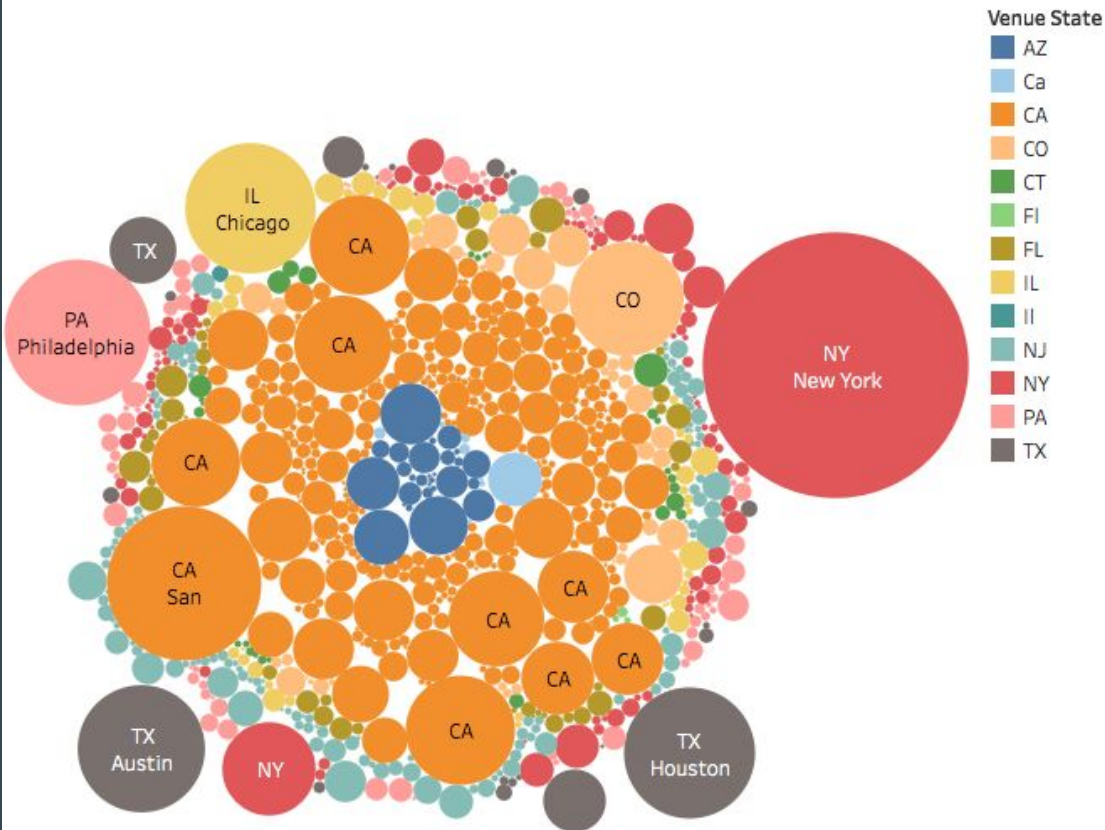
```
sqlContext.sql("CREATE TABLE groupCity AS SELECT group.name, venue.city FROM events_temp")
```


Sample Data Table - Visualization via Tableau



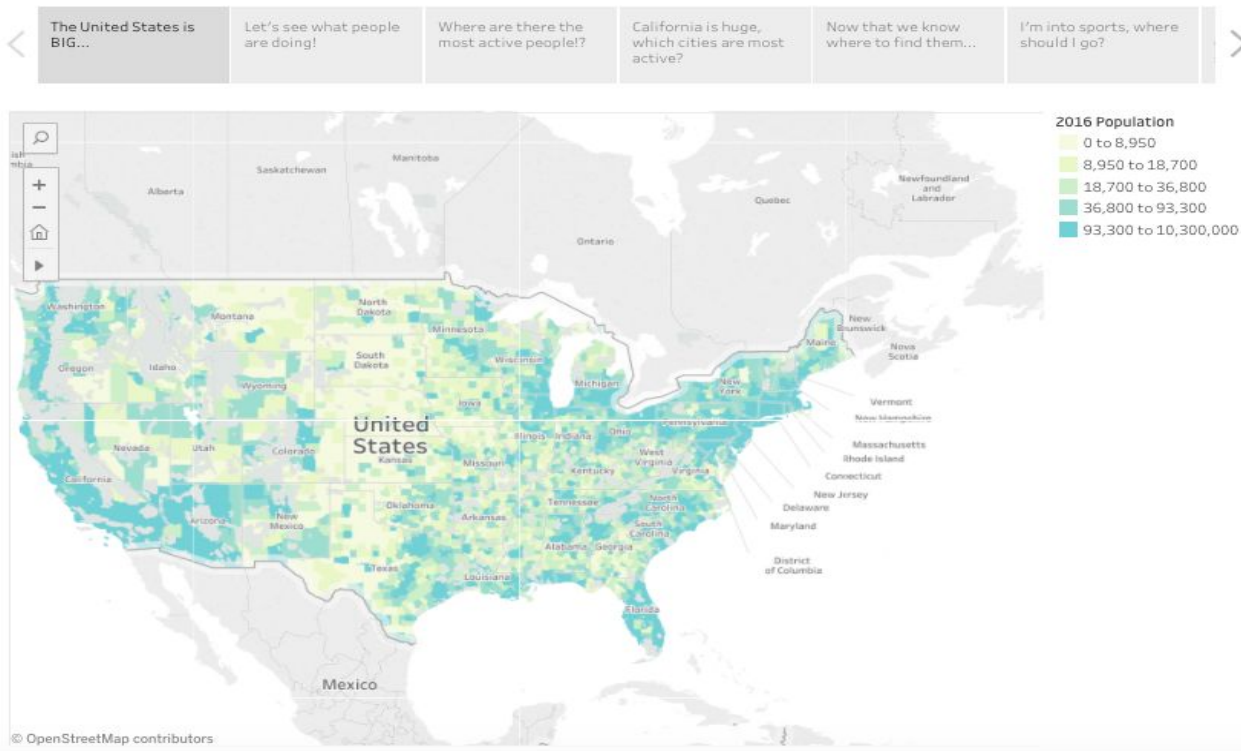
Sample Data Viz - Tableau

Most Popular Cities and States



Output: Stories

Find a New Location



Output: Dashboards

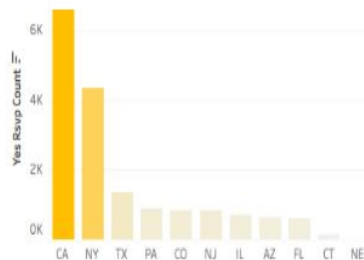
Find Your Niche



Choose your category

(All)

Top States for Category



Most Popular Days



Most Popular Meetup Times



with Love

What do you like to do?

(All)

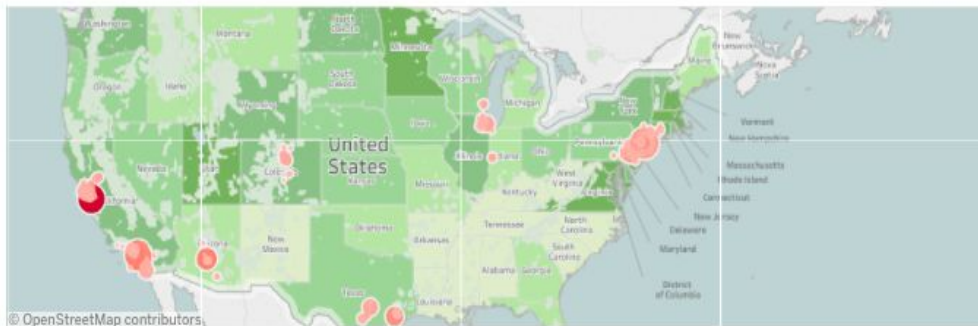
Which states would you live in?

(Multiple values)

What days are you free?

(All)

Best States for Matches



Complexity and Storage Needs

- Complexity
 - Project complexity is moderated by design
 - Spouts are well behaved, present well formed data
 - Data velocity is throttled by API throttling
 - System is designed to handle higher velocity
 - Database kept co-located with Ingest, provides managed latency
 - Overcomes inability to extract data over relatively slow internet
 - Hive and Map Reduce well understood but not pushed to limit
 - Hive latency issues due to the map reduce framework not magnified
- Storage
 - Tested at order of 86Mb/hour, can be scaled by 100x
 - To meet higher scale, Ingest layer may need distributed Storage architecture
 - Distributed storage provides capability to handling diverse data sets
 - Can push magnification of I/O rate by 100x

Challenges and Remedies

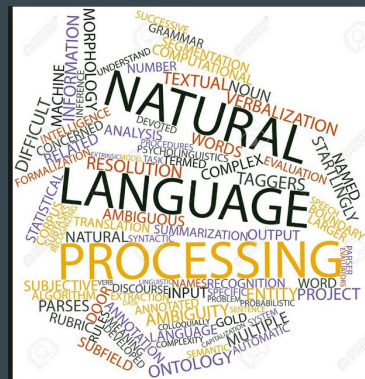
- API
 - Meetup api is limited in scope. For instance, querying by epoch is not supported!
 - Some required fields needs special parameter tags and are expensive to extract
 - API requests are limited to 200 per hour, need to throttle rate
- Processing
 - Multiple points of interface: PyMongo/ MongoHadoop/ PyMongoSpark/ Spark interface, many points of data exchange
- Front End
 - Despite the processing, still have minor issues with redundancy
 - Solution - opt for stricter constraints
 - Timezones
 - Known problem with all programming
 - ST solution- convert time from epoch to date-time then adjust by state
 - LT Solution- create tables for offsets based upon zipcodes and recognition of

Future Product Improvements

- Extend filtering and grouping algorithm for Precise Calculations
- Support Deeper Insights via:
 - Group analysis
 - Keyword Search
 - Domain knowledge
- Seek out and group Demographic Information
 - Focus on individuals, not groups
 - Pull information from profiles and connect to attended events
- Refine Ease of Use
 - Even with more robust information, we have a business tool
 - Need to make user friendly with controlled/verified output
- Extend application capability
 - Extract public sentiment

Scaling: Now and Future

- Support multiple data spouts, each with own schema
- Extend I/O Scaling: Obtaining and saving data
 - Scale the Database architecture, possibly distributed
 - Storage buffer in Ingest layer
- API
 - Extend API calls for more subfield level queries
- Analysis
 - Extend to predictive and prescriptive conclusions
- Extend towards capability for custom output
 - GUI output for casual applications
 - CDN for advanced users



Did Project meet set goals?

- ETL infrastructure for Meetup data
- Summarization of data across 12 large metropolitan areas
- Summarization of recent data
 - What activities are communities throughout the United States excited about?
 - Most prominent meetings in the US in December 2016
 - Does how people connect through activity differ by region?

