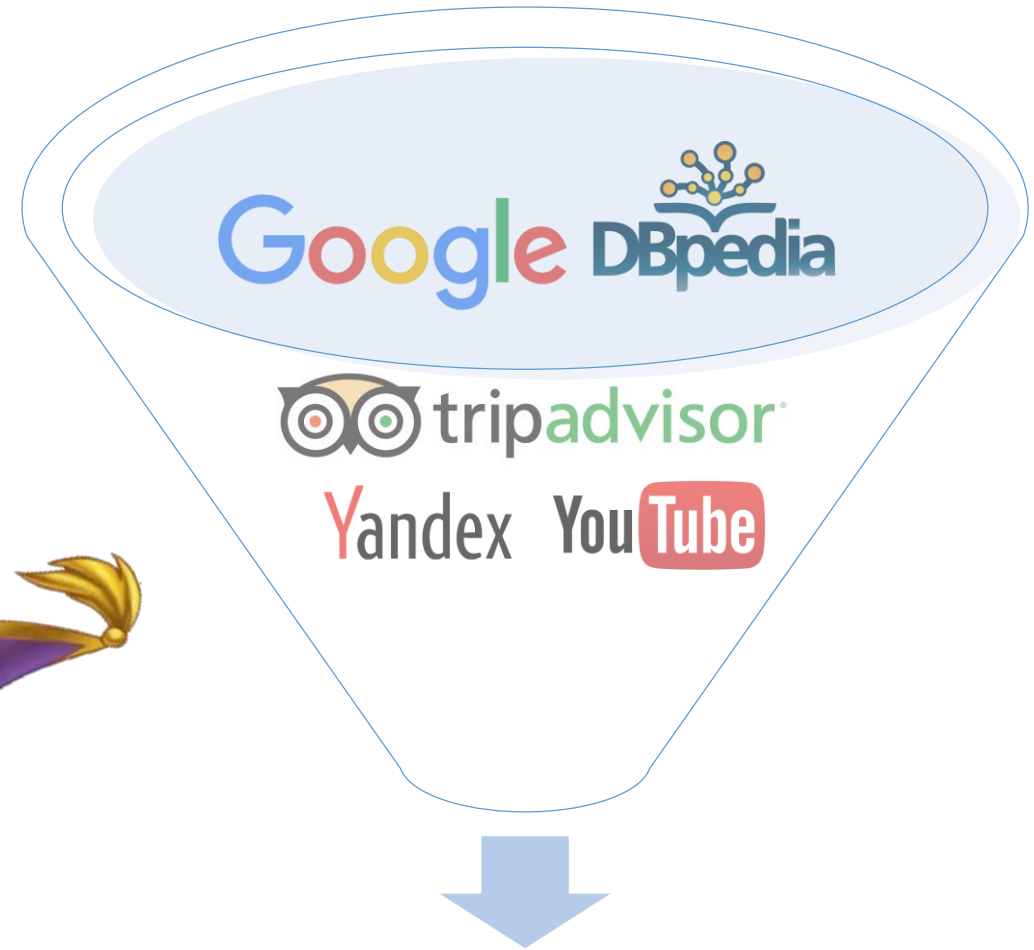




Rashmi Nalwad  
Ravi Raju Krishna



## Enroute-Genie

**A Coherent 1-Stop solution to explore the world**

[www.enroute-genie.com](http://www.enroute-genie.com)  
<https://github.com/raviraju/Enroute-Genie>



Aladdin's ancient carpet was found  
in the Cave of Wonders' treasure  
room

Our Software Carpet emerged  
from IIW course at WPH B27







## Enroute Genie Mashup

A Data Integration Solution powered by     

**Start From:**  **End At:**



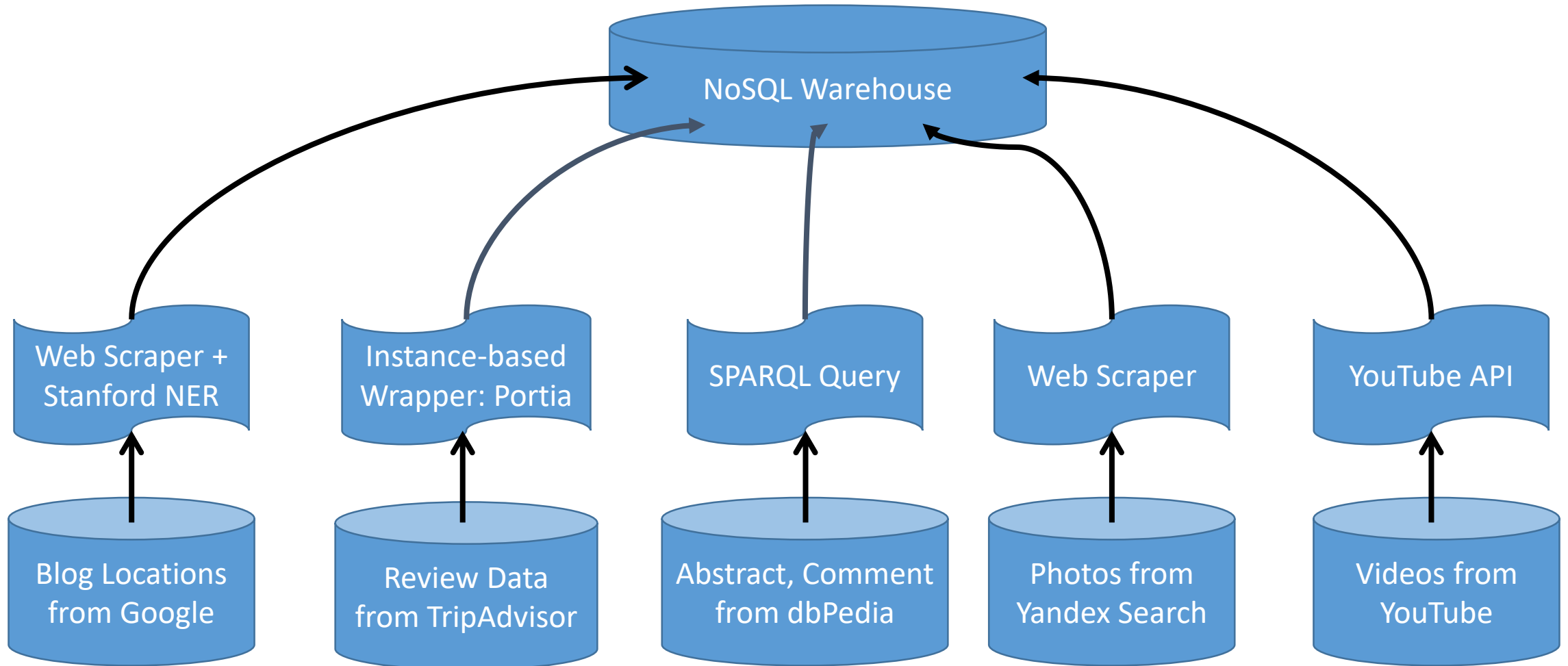
Learning Objectives of Course	Project Objectives
Understand the theory and techniques of traditional data integration, including logical view integration, schema mapping, and <b>record linkage</b>	Link Records from Google, TripAdvisor & DBPedia
Understand the foundations and techniques of the <b>Semantic Web</b> , including RDF, OWL, <b>SPARQL</b> , <b>linked data</b>	Query data from DBPedia – the semantic web version of Wikipedia, and add HTTP URI to rdf sources adhering to principles of linked data
Understand the theory and application of the state-of-the-art software and tools for <b>information extraction</b>	Scrape Location mention from blogs using customized NER
Understand the algorithms and techniques for <b>data cleaning</b> , source modeling, semi-structured extraction, and information extraction	Configure an Instance based Wrapper Induction Extraction system to fetch structured information from TripAdvisor and use Trifacta Wrangler to refine the results.
For <b>any given integration problem</b> , be able to select and apply the most relevant information integration techniques to solve that problem	Build Enroute-Genie

# Modest Goals of data integration system

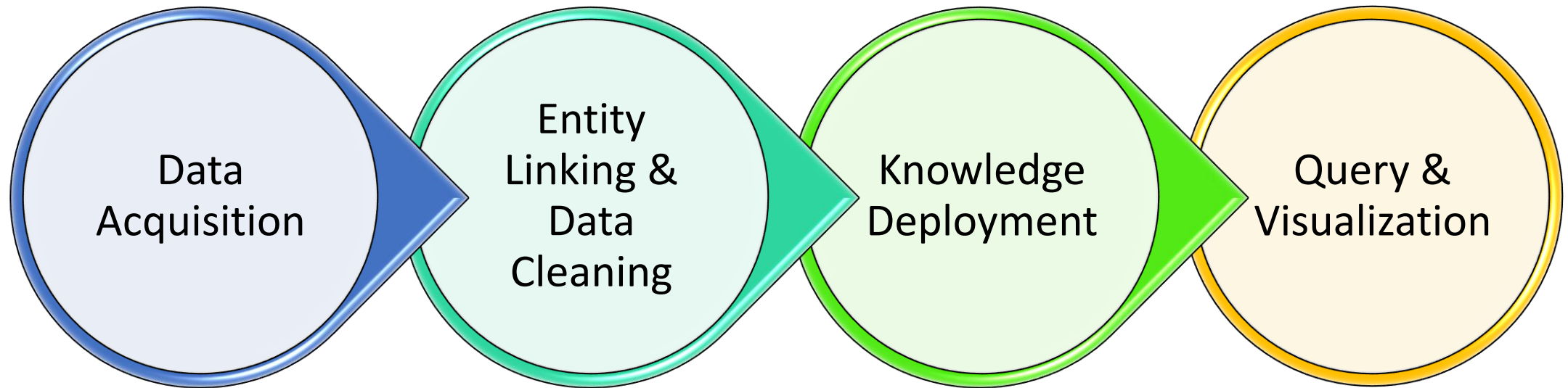
- User Effort : *To build tools that reduce the effort required to integrate a set of data sources. Make it easy to add new sources*
    - \* Currently addressing State of California.
    - \* Can be easily extended to other states.
    - \* Triggered extension to Karnataka, India, during development.
  - Accuracy : *To improve the ability of the system to answer queries in uncertain environments.*
    - \* We are using current state of art for NER – StanfordCoreNLP to identify locations.
    - \* Our accuracy improves adapting to future advances in NER
- [NOTE : See Evaluation Results for more info...]



# Architecture



# Pipeline





# Data Sources

## **Variety** dimension in 3 V's of Big Data

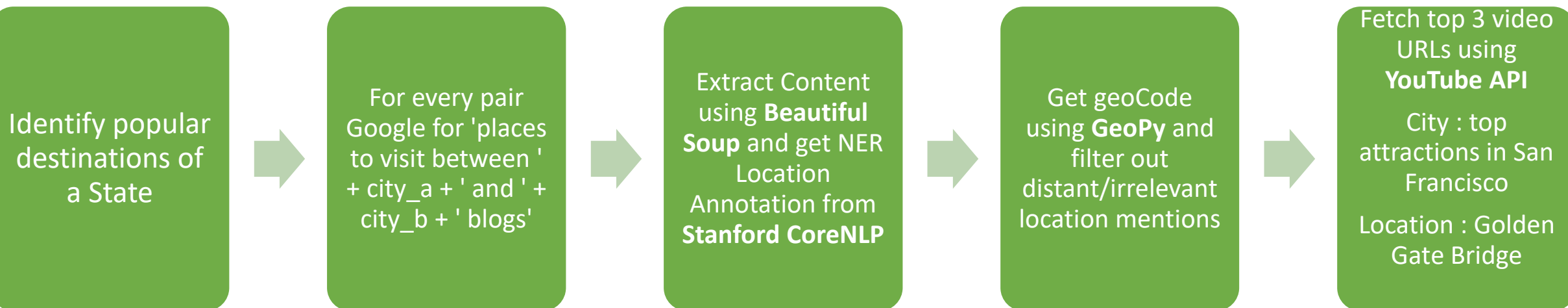
1. Unstructured : Blog text Scraped from Google search results
2. Structured : Abstract and Comment of a city from dbPedia, TripAdvisor Review Data
3. Multimedia : Photos and YouTube Videos

Accessed:

## Information Extraction

1. Scraper(Blog Links from Google & ImageUrls from Yandex)
2. Mine the Deep Web Data : Wrapper(TripAdvisor data)
3. SPARQL Query(dbPedia Info)
4. Web Service API(Youtube API)

# Data Acquisition



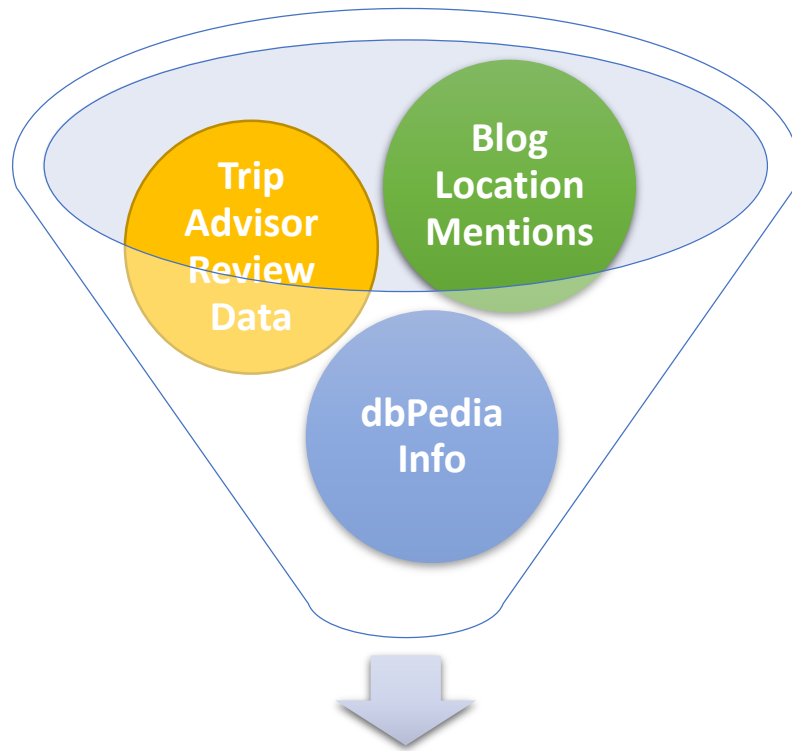
Configure an Instance-Based Data extraction wrapper-**Portia** to extract attraction review metadata



Clean Data using **Trifacta Wrangler**

Query dbPedia for Abstract and Comment of locations partOf California Using **SPARQL**

# Entity Linking – using FRIL



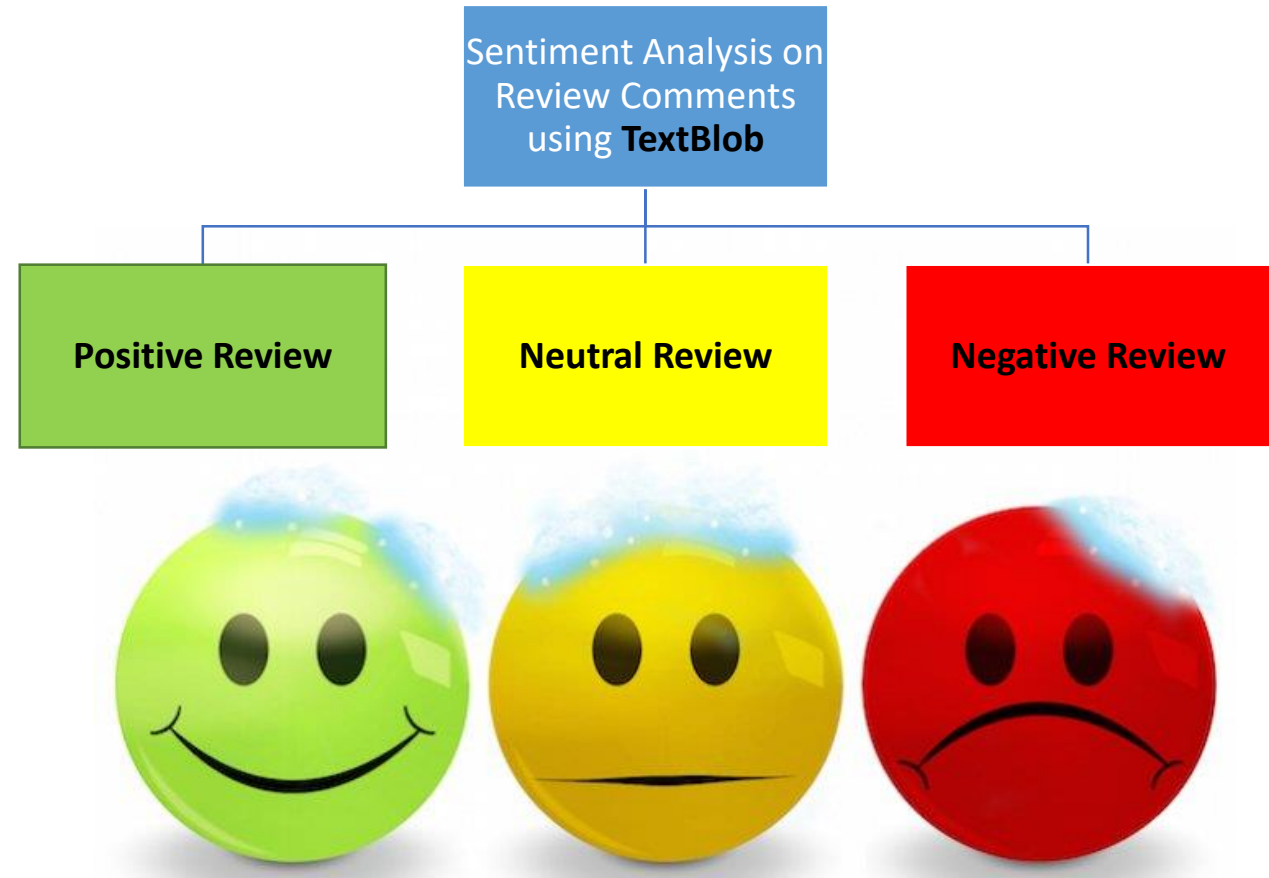
Record Linked Data

**Edit-Distance :** (mention\_name, attraction\_name),  
(match-level-start=0.1, math-level-end=0.3)  
weight="50"

**Numeric-Distance :** (mention\_latitude, attraction\_latitude)  
exactMatch  
weight="25"

**Numeric-Distance :** (mention\_longitude, attraction\_longitude)  
exactMatch  
weight="25"

# Enrich Record Linked Data





# Knowledge Deployment

- Integrated Data loaded to NoSQL-CouchDB
- Web Service to share our integrated data with other applications:

HTTP API Access to our DB

- [https://enroutegenie.cloudant.com/enroutegenie/all\\_docs](https://enroutegenie.cloudant.com/enroutegenie/all_docs)
- [https://enroutegenie.cloudant.com/enroutegenie/los%20angeles and san%20francisco](https://enroutegenie.cloudant.com/enroutegenie/los%20angeles%20and%20san%20francisco)

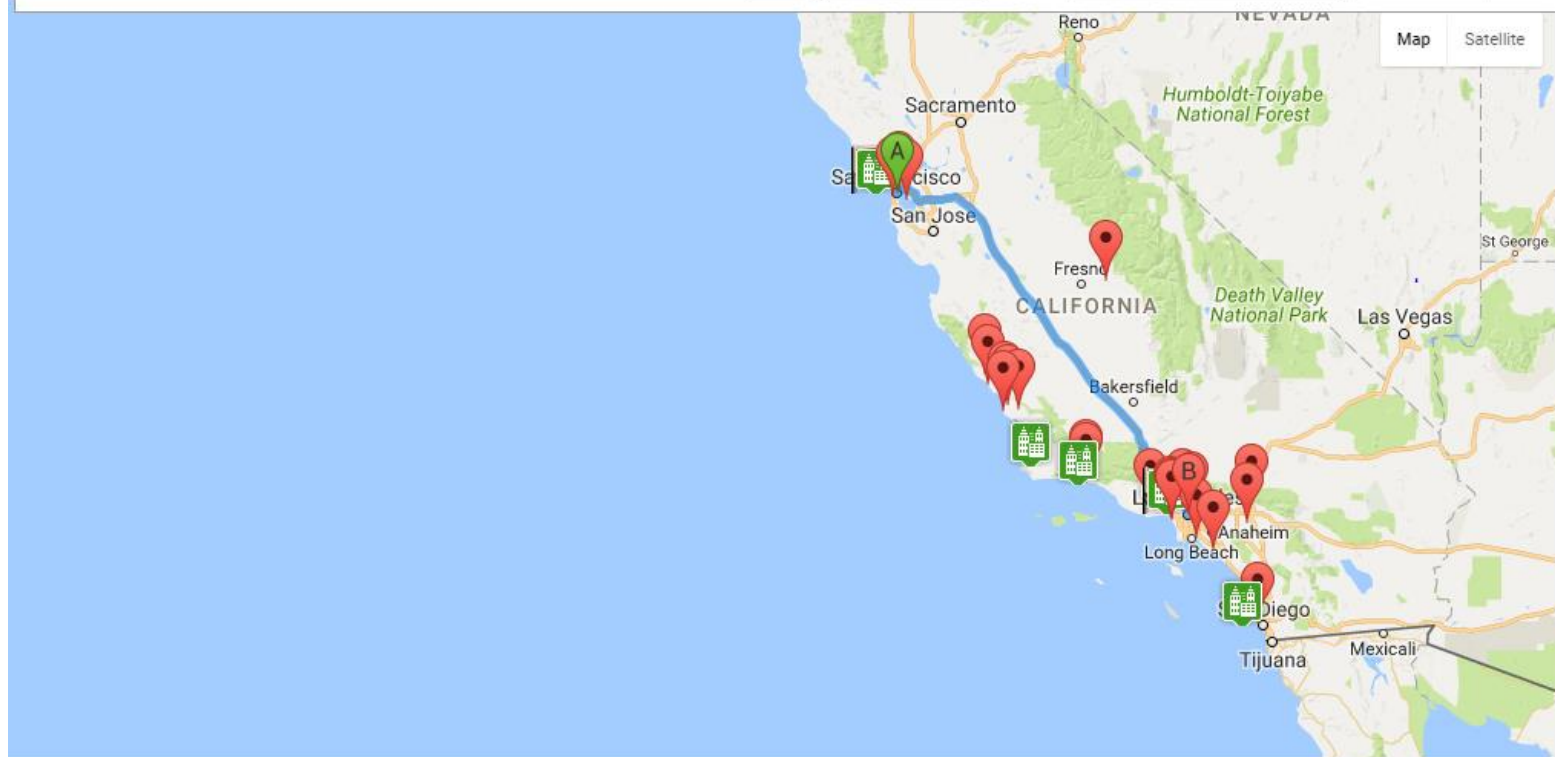


# Query & Visualization

Enroute Genie Mashup

A Data Integration Solution powered by Google, TripAdvisor, DBpedia, Yandex, YouTube

Start From: los angeles End At: san francisco Go Genie !!!



Map Satellite

1 Dr Carlton B Goodlett Pl, San Francisco, CA 94102, USA

382 mi. About 5 hours 39 mins

1. Head south on Polk St 0.2 mi toward Grove St
2. Continue onto 10th St 0.6 mi
3. Turn left onto Bryant St 0.2 mi
4. Turn left onto the 0.2 mi Interstate 80 E ramp to Oakland
5. Merge onto I-80 E 6.8 mi
6. Take the Interstate 580 0.9 mi E exit toward CA-24/Hayward/Stockton

# City Attractions

Enroute Genie Mashup

A Data Integration Solution powered by Google tripadvisor DBpedia Yandex YouTube

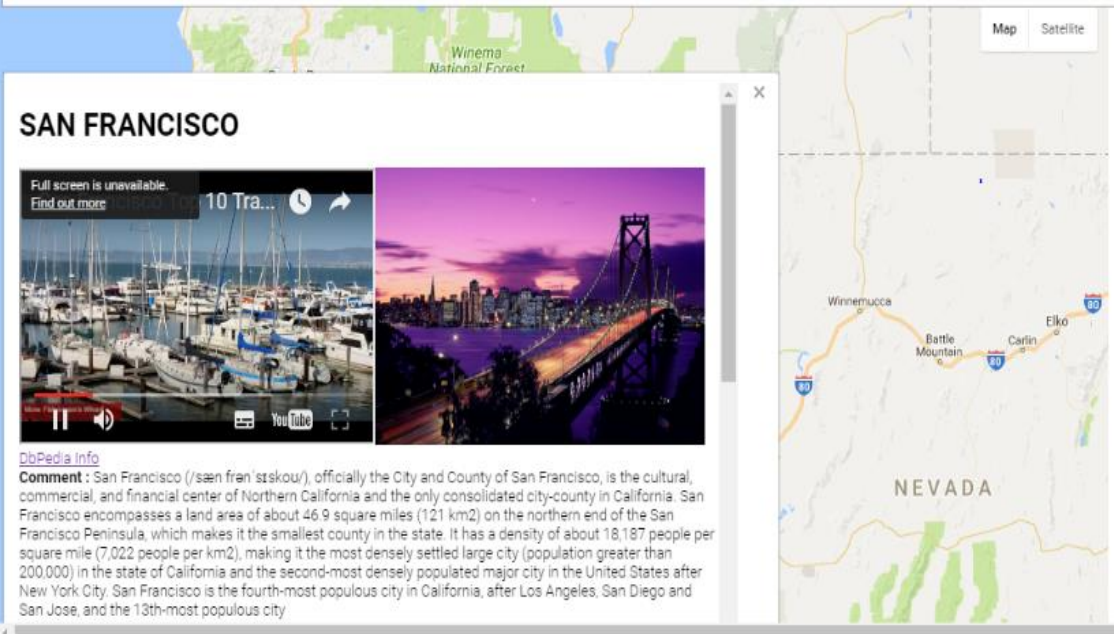
Start From: los angeles End At: san francisco Go Genie !!!

**SAN FRANCISCO**

Full screen is unavailable. Find out more 10 Tra...

**Dbpedia info**

**Comment :** San Francisco (/sæn frənˈsɪskəʊ/), officially the City and County of San Francisco, is the cultural, commercial, and financial center of Northern California and the only consolidated city-county in California. San Francisco encompasses a land area of about 46.9 square miles (121 km<sup>2</sup>) on the northern end of the San Francisco Peninsula, which makes it the smallest county in the state. It has a density of about 18,187 people per square mile (7,022 people per km<sup>2</sup>), making it the most densely settled large city (population greater than 200,000) in the state of California and the second-most densely populated major city in the United States after New York City. San Francisco is the fourth-most populous city in California, after Los Angeles, San Diego and San Jose, and the 13th-most populous city



Enroute Genie x D About: San Francisco x

dbpedia.org/page/San\_Francisco

DBpedia Browse using Formats Faceted Browser Sparql Endpoint

## About: San Francisco

An Entity of Type : Consolidated city-county, from Named Graph : <http://dbpedia.org>, within Data Space : [dbpedia.org](http://dbpedia.org)

San Francisco (/sæn frənˈsɪskəʊ/), officially the City and County of San Francisco, is the cultural, commercial, and financial center of Northern California and the only consolidated city-county in California. San Francisco encompasses a land area of about 46.9 square miles (121 km<sup>2</sup>) on the northern end of the San Francisco Peninsula, which makes it the smallest county in the state. It has a density of about 18,187 people per square mile (7,022 people per km<sup>2</sup>), making it the most densely settled large city (population greater than 200,000) in the state of California and the second-most densely populated major city in the United States after New York City. San Francisco is the fourth-most populous city in California, after Los Angeles, San Diego and San Jose, and the 13th-most populous city

Property	Value
<a href="#">dbo:PopulatedPlace/areaMetro</a>	▪ 9128.1540960682
<a href="#">dbo:PopulatedPlace/areaTotal</a>	▪ 600.592342905815
<a href="#">dbo:PopulatedPlace/populationDensity</a>	▪ 7022.039957411464
<a href="#">dbo:abstract</a>	▪ San Francisco (/sæn frənˈsɪskəʊ/), officially the City and County of San Francisco, is the cultural, commercial, and financial center of Northern California and the only consolidated city-county in California. San Francisco encompasses a land area of about 46.9 square miles (121 km <sup>2</sup> ) on the northern end of the San Francisco Peninsula, which makes it the smallest county in the state. It has a



# Location Mentions

Enroute Genie Mashup

A Data Integration Solution powered by Google, TripAdvisor, DBpedia, Yandex, YouTube

Start From: los angeles End At: san francisco Go Genie !!!

## MORRO BAY STATE PARK

Full screen is unavailable. Find out more

Morro Bay State Park A Campground Fav!

Travel Small-Live Big!

Known For "Specialty Museums , Museums"

TripAdvisor Rank 4

TripAdvisor No of Reviews 315

No of Blogs mentioned 5

Contact Info +1 805-772-2694


Located In Morro Bay

**A Positive Review** Camping, hiking, golf and relaxation in the midst of an eucalyptus grove next to beautiful Morro Bay with beautiful views of the San Luis Obispo County hills that stay green most... [read more](#)

**A Neutral Review** Lots of open areas to hike in park without paying for parking. [read more](#)

**Address** Morro Bay State Park, South Bay Boulevard, Baywood Park, San Luis Obispo County, California, 93442, United States of America

**Location Type** national\_park



Enroute Genie Mashup

WINNER WONDERLAND Win Big When You Book on TripAdvisor Open Your

tripadvisor Morro Bay State Park: Address, Phone Number, Specialty Museum Reviews

Morro Bay Hotels Flights Vacation Rentals Restaurants Things to Do Forum Best of 2016 More

Find: Things to Do Near: Morro Bay, California

United States > California (CA) > San Luis Obispo County > Morro Bay > Things to Do in Morro Bay > Morro Bay State Park


## Morro Bay State Park

Is this your business?

316 Reviews #4 of 35 things to do in Morro Bay Certificate of Excellence

Specialty Museums, Museums

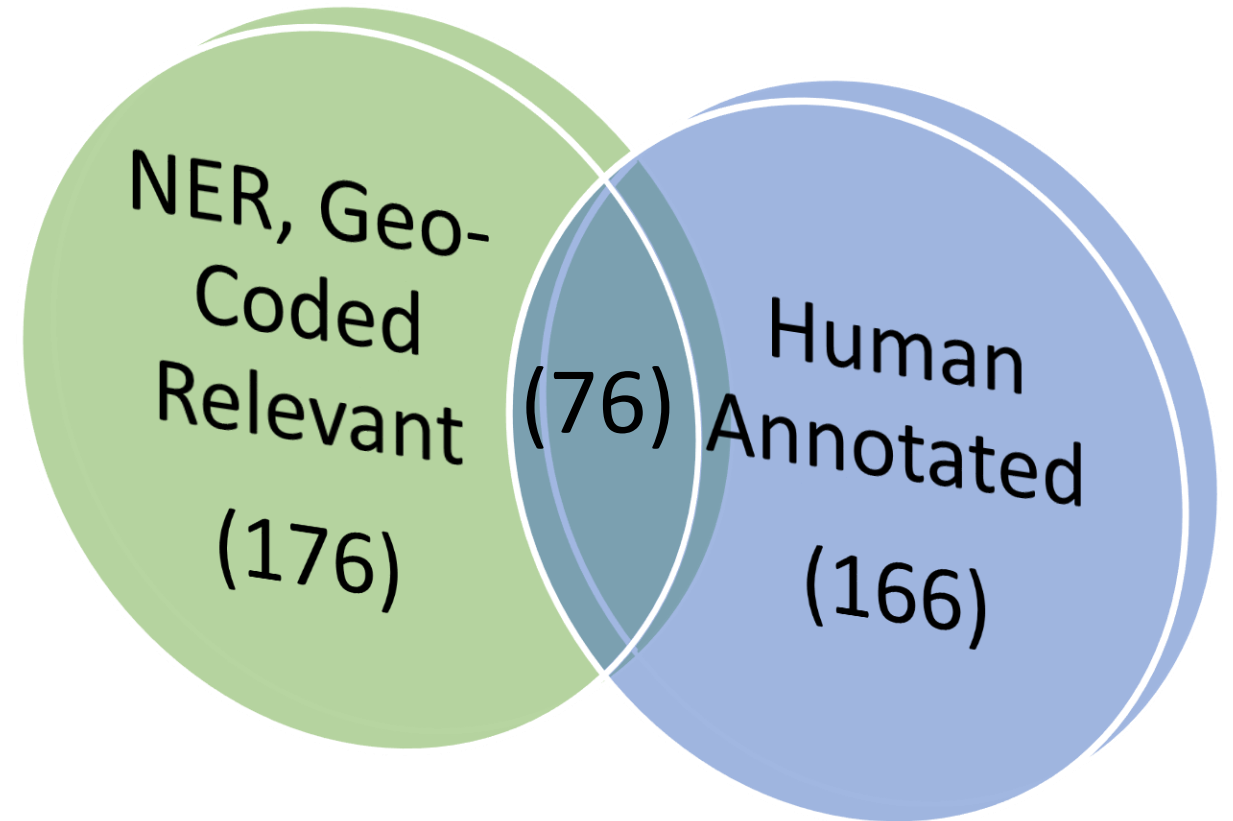
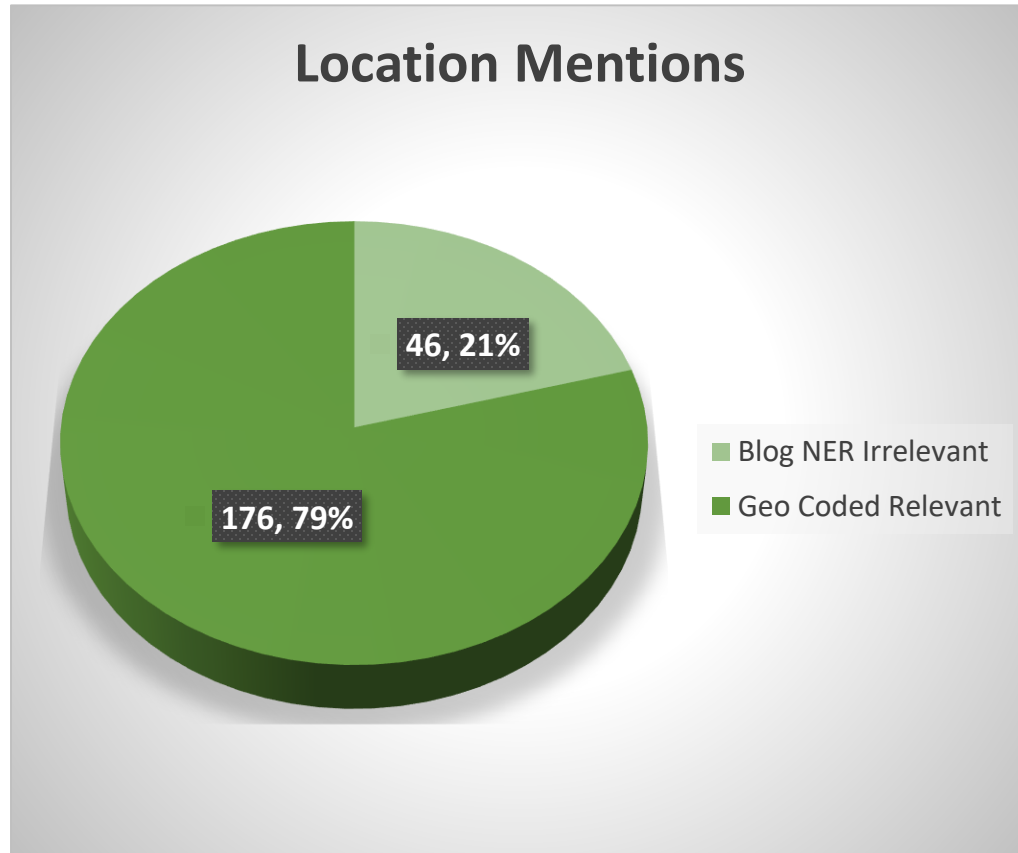
Overview Reviews (316) Q&A (7) Location



All visitor photos (97)



# Evaluation – Los Angeles and San Francisco



# Data Integration or Tools Integration ?

Feature	Library	Reference
Blog Text	Google Scraper-"Bing Search Engine"	<a href="https://github.com/NikolaiT/GoogleScraper">https://github.com/NikolaiT/GoogleScraper</a> <a href="https://www.crummy.com/software/BeautifulSoup/bs4/doc/">https://www.crummy.com/software/BeautifulSoup/bs4/doc/</a>
	Beautiful Soup	
NER Location	Stanford CoreNLP	<a href="http://stanfordnlp.github.io/CoreNLP/">http://stanfordnlp.github.io/CoreNLP/</a>
Geo Coding	GeoPy	<a href="http://geopy.readthedocs.io/en/latest/">http://geopy.readthedocs.io/en/latest/</a>
TripAdvisor Data	Portia, Scraping Hub	<a href="https://doc.scrapinghub.com/">https://doc.scrapinghub.com/</a> <a href="https://doc.scrapinghub.com/portia.html">https://doc.scrapinghub.com/portia.html</a>
Data Cleaning	Trifacta Wrangler	<a href="https://www.trifacta.com/start-wrangling/">https://www.trifacta.com/start-wrangling/</a>
Record Linking	FRIL	<a href="http://fril.sourceforge.net/">http://fril.sourceforge.net/</a>
Sentiment Analysis	Text Blob	<a href="https://textblob.readthedocs.io/en/dev/api_reference.html#textblob.blob.TextBlob.sentiment">https://textblob.readthedocs.io/en/dev/api_reference.html#textblob.blob.TextBlob.sentiment</a>
Videos	YouTube API	<a href="https://developers.google.com/youtube/v3/code_samples/python#search_by_keyword">https://developers.google.com/youtube/v3/code_samples/python#search_by_keyword</a>
Images	Google Scraper-"Yandex Search Engine"	<a href="https://github.com/NikolaiT/GoogleScraper">https://github.com/NikolaiT/GoogleScraper</a>



# More on Grading

- This is a hard class, but you will learn a lot!
  - Principles and theory
    - Technical readings and lectures (quizzes, final exam)
  - Putting principles into practice
    - Homeworks and project!
- Grade distribution

94 - 100 = A	74 - 76 = C
90 - 93 = A-	70 - 73 = C-
87 - 89 = B+	67 - 69 = D+
84 - 86 = B	64 - 66 = D
83 - 83 = B-	60 - 63 = D-
77 - 79 = C+	Below 60 is an F

Questions ?

