



# Yelper

## A Collaborative Filtering Based Recommendation System

Team DeepBlue

Chuan Sun

Capstone Project @ NYC Data Science Academy

9/20/2016



# Agenda

## What?

What are the components for  
Yelper?

## Why?

Why we need  
recommendation?

## Demo

User-business network; Yelper  
main page; Simulation of  
users' recommendation  
requests handling

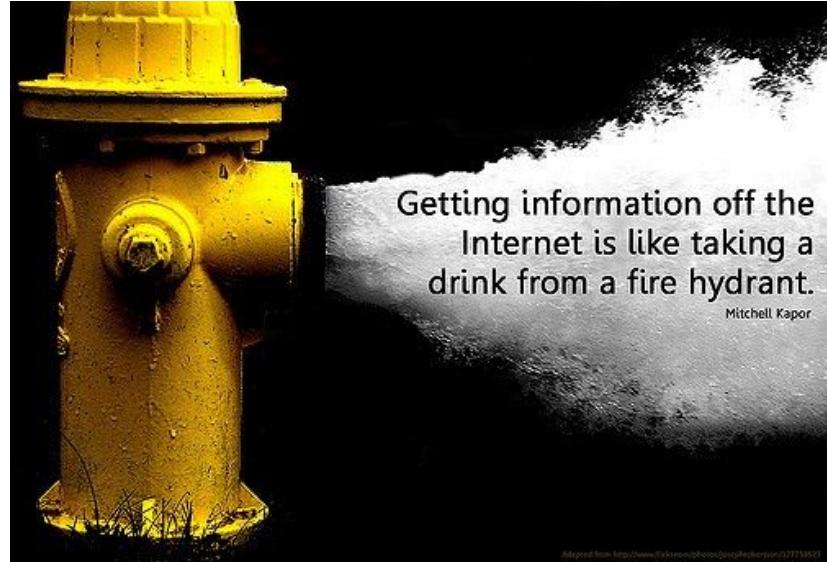
## How?

How to build Yelper?

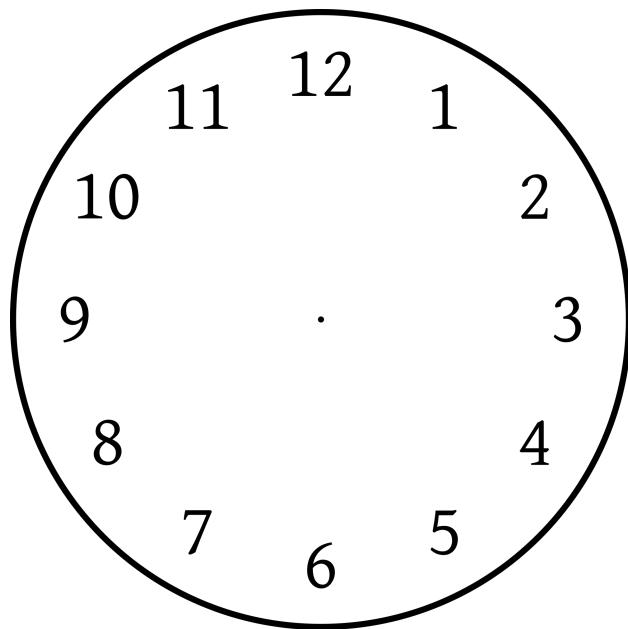
## Summary & Future works

What we learn and what's  
next?

# Information overload is a real phenomenon which prevents us from taking decisions or actions



But life is short. We only have 24 hours per day

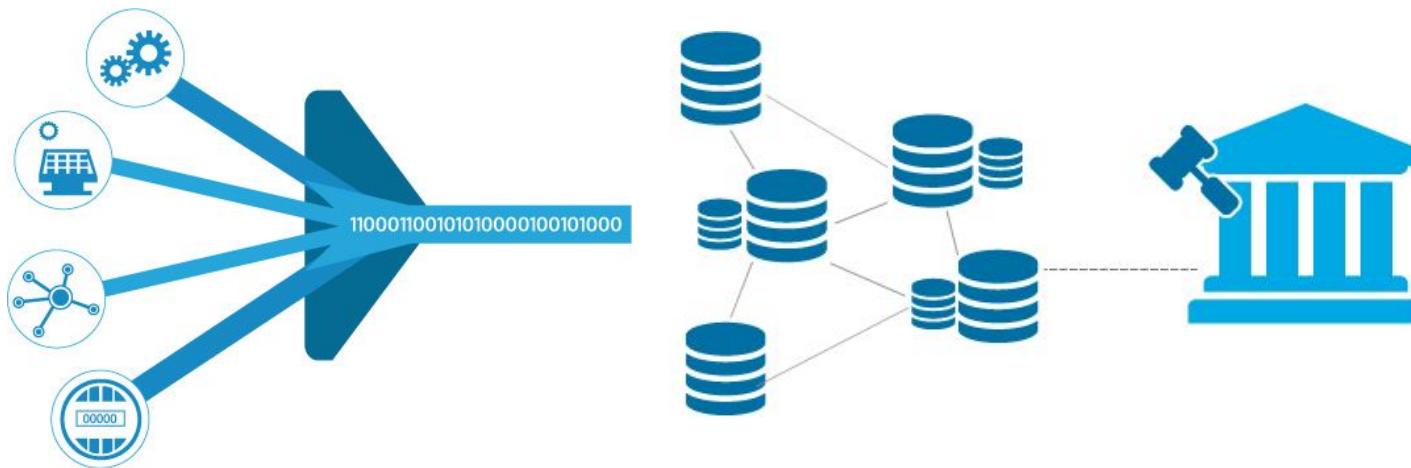


Recommendation is becoming extremely common in recent years



# Real-time recommendation system may become a new normal in data era

- Gain real-time insight
- Enable rapid development
- Perform real-time analytics



# Agenda

## What?

What are the components for  
Yelper?

## Why?

Why we need  
recommendation?

## Demo

User-business network; Yelper  
main page; Simulation of  
users' recommendation  
requests handling

## How?

How to build Yelper?

## Summary & Future works

What we learn and what's  
next?

# Yelp Challenge 2016 Dataset is a good sandbox to build recommendation engine

## The Challenge Dataset:

- **2.7M** reviews and **649K** tips by **687K** users for **86K** businesses
- **566K** business attributes, e.g., hours, parking availability, ambience.
- Social network of **687K** users for a total of **4.2M** social edges.
- Aggregated check-ins over time for each of the **86K** businesses
- **200,000** pictures from the included businesses

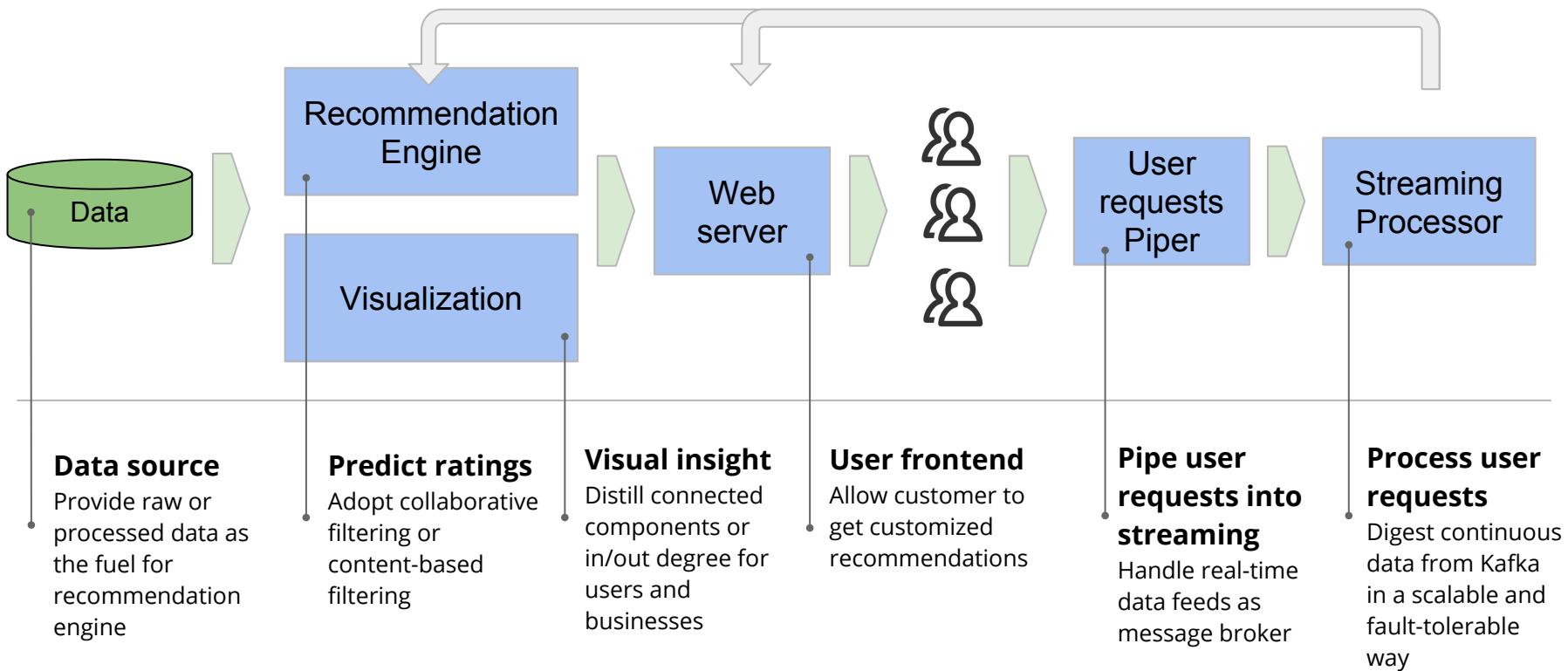
## Cities:

- U.K.: Edinburgh
- Germany: Karlsruhe
- Canada: Montreal and Waterloo
- U.S.: Pittsburgh, Charlotte, Urbana-Champaign, Phoenix, Las Vegas, Madison

## Best of Yelp: New York

Food	Food	See More
 Food	 Food	<a href="#">See More</a>
 Nightlife	 1. <a href="#">Coffee Project New York</a>	
 Restaurants	 ★★★★☆ 213 reviews	
 Shopping	 Or if you're feeling adventurous as i did, try their deconstructed latte.	
 Active Life	 2. <a href="#">Levain Bakery</a>	
 Arts & Entertainment	 ★★★★★ 4622 reviews	
 Automotive	 Crispy on the outside, served warm with a moist crumbly inside.	
 Beauty & Spas	 3. <a href="#">Borgatti's Ravioli &amp; Egg Noodles</a>	
 Education	 ★★★★☆ 109 reviews	
 Event Planning & Se...	 ONLY FEW pounds of the fresh pasta here in Borgatti's Ravioli & Egg Noodles!	
 Health & Medical	 4. <a href="#">Foods of NY Tours</a>	
 Home Services	 ★★★★☆ 203 reviews	
 Local Services	 Darrell was a fantastic guide and we loved Chelsea.	
 More Categories	 5. <a href="#">Coney Shack</a>	
	 ★★★★☆ 228 reviews	
	 Faves: the Vietnamese short rib taco and beer battered fish taco.	

# Yelper consists of 5 major components



# Agenda

## What?

What are the components for  
Yelper?

## Why?

Why we need  
recommendation?

## Demo

User-business network; Yelper  
main page; Simulation of  
users' recommendation  
requests handling

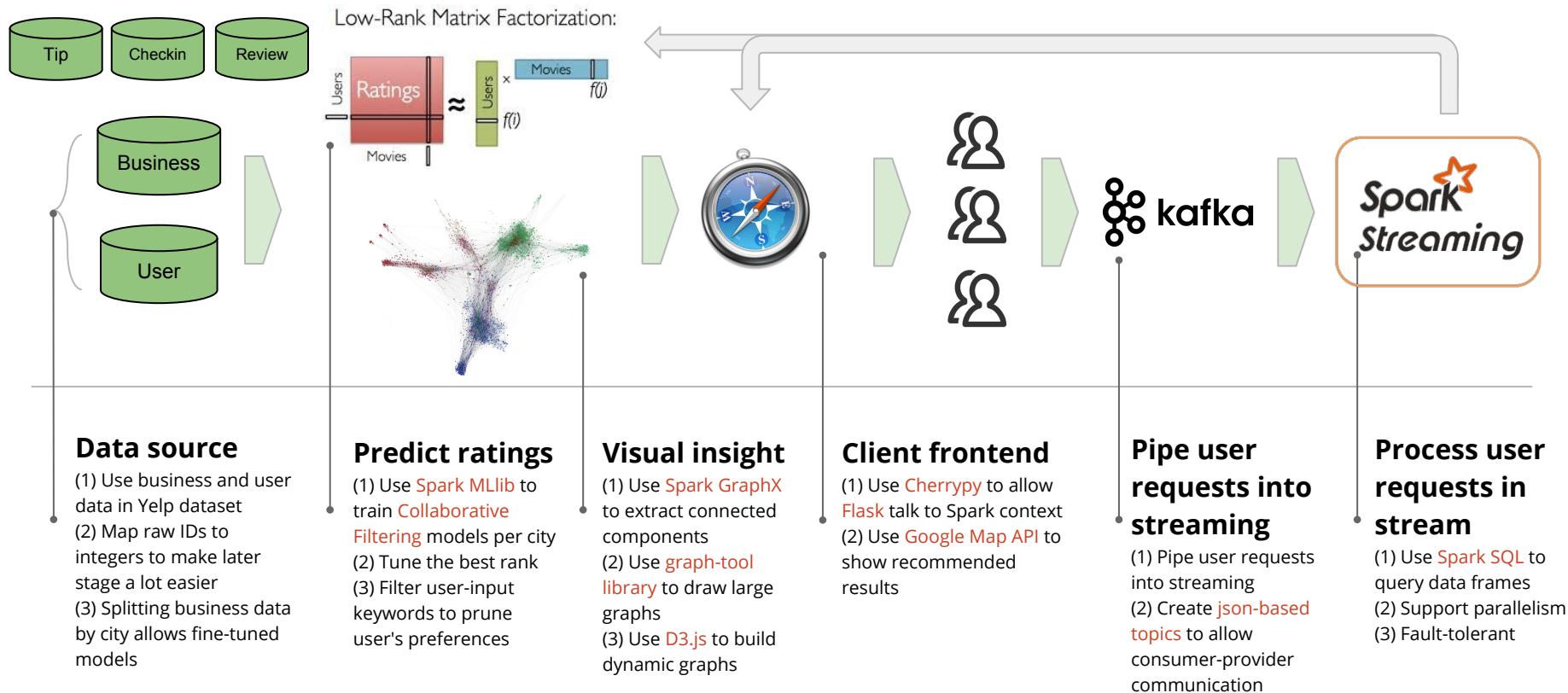
## How?

How to build Yelper?

## Summary & Future works

What we learn and what's  
next?

# Yelper is built on top of Spark platform



# Agenda

## What?

What are the components for  
Yelper?

## Demo

User-business network; Yelper  
main page; Simulation of  
users' recommendation  
requests handling

## Why?

Why we need  
recommendation?

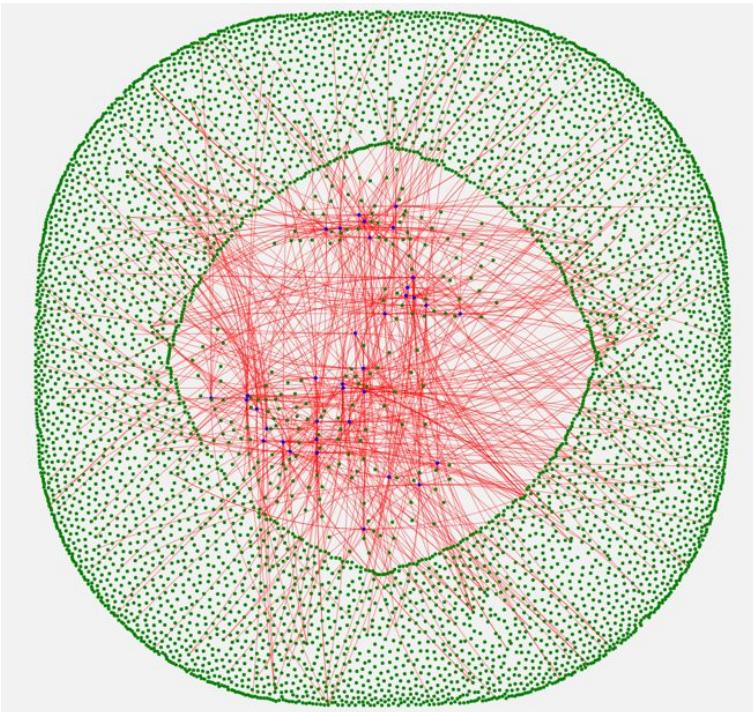
## How?

How to build Yelper?

## Summary & Future works

What we learn and what's  
next?

# The city-wise user-business network in Yelp tells a lot about a city



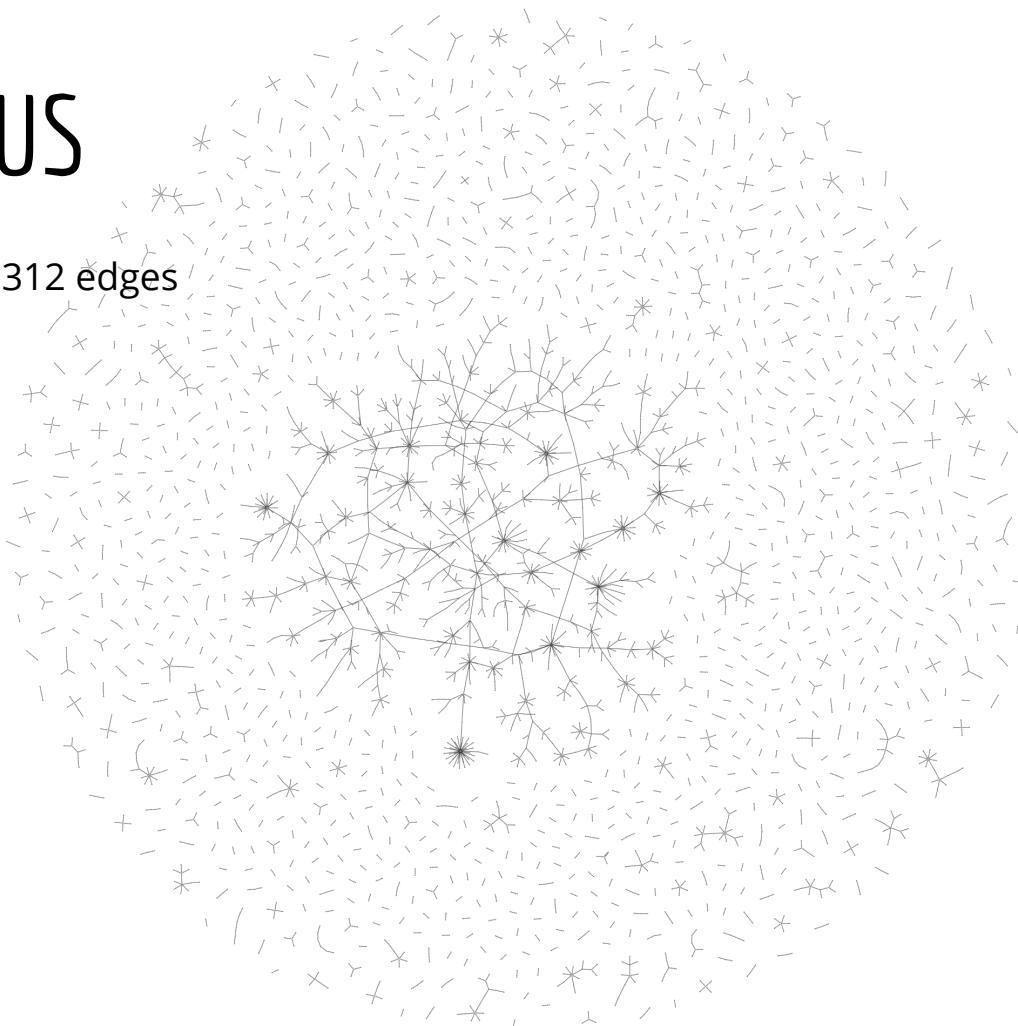
- Motivation
  - User-business interaction may reflect the economy trend in a city
- Each city is distinct, so does its user-business network
- Nodes:  $u_1, u_2, \dots, b_1, b_2, \dots$
- Edges:  $u_1 \rightarrow b_1, u_1 \rightarrow b_2, \dots$
- The network is a bipartite
- If Yelp provides the timestamp of rating data, we may see the evolution of economy development

Demo 1: User-business network visualization in the city of Madison, US

# Charlotte, US

Randomly selected 3312 edges

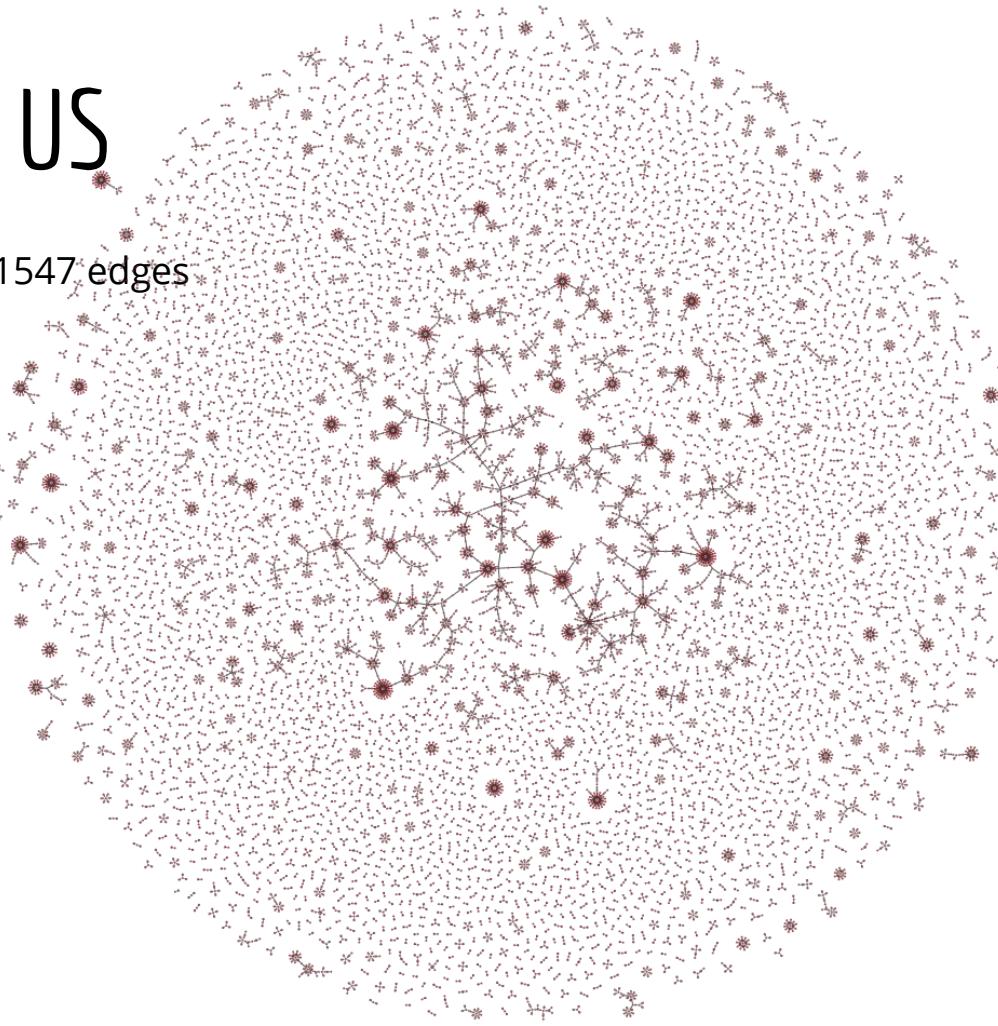
(1% of total ratings)



# Las Vegas, US

Randomly selected 11547 edges

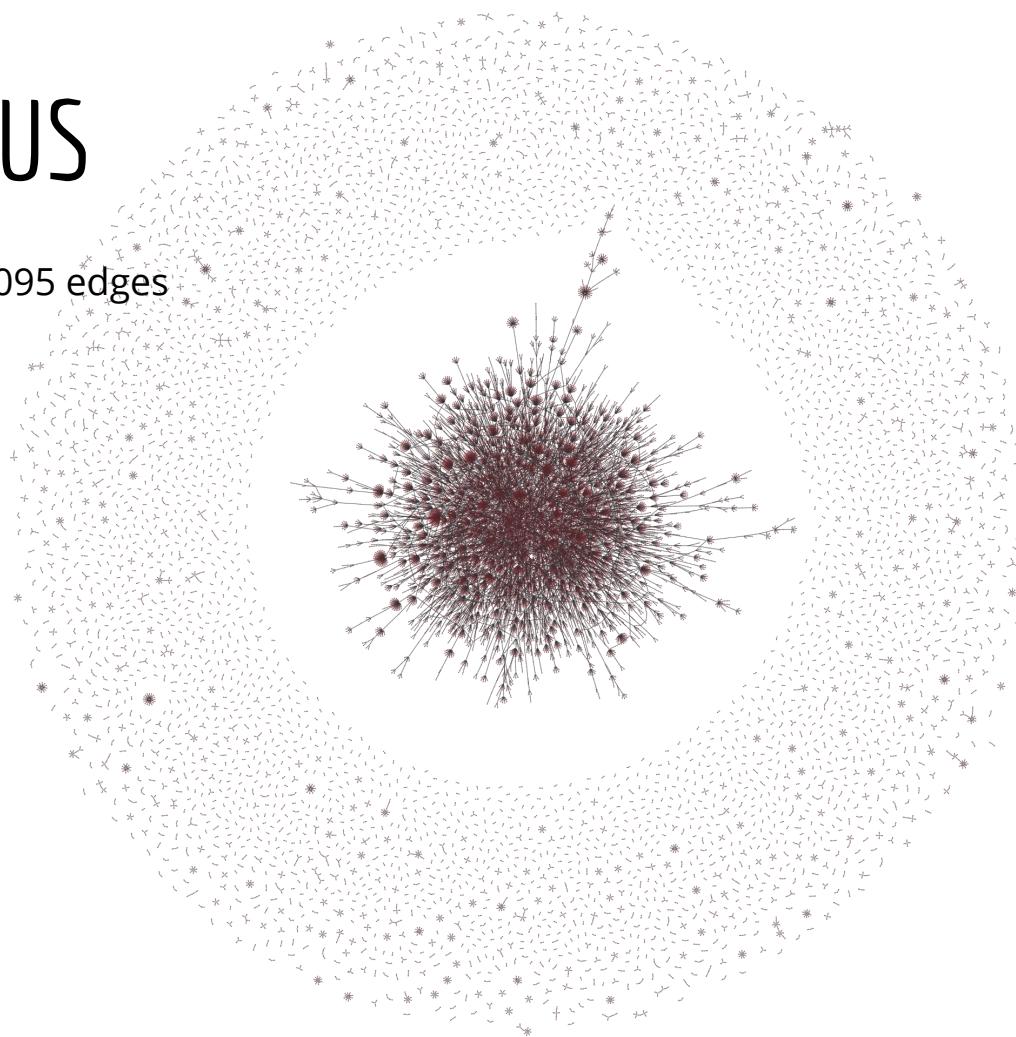
(1% of total ratings)



# Las Vegas, US

Randomly selected 23095 edges

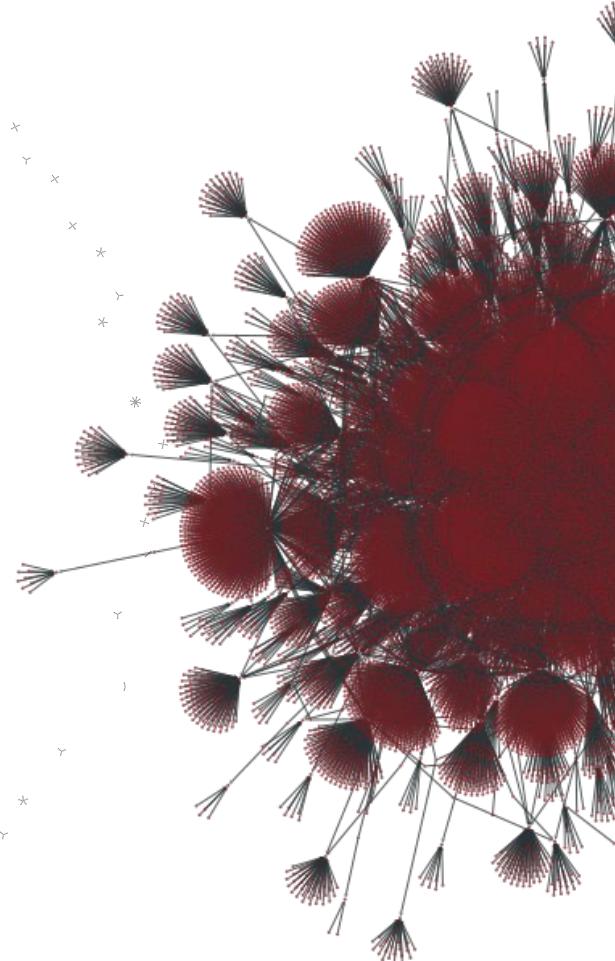
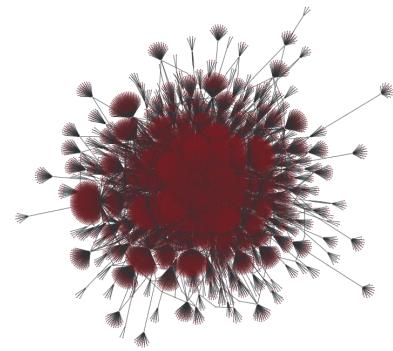
(2% of total ratings)



# Las Vegas, US

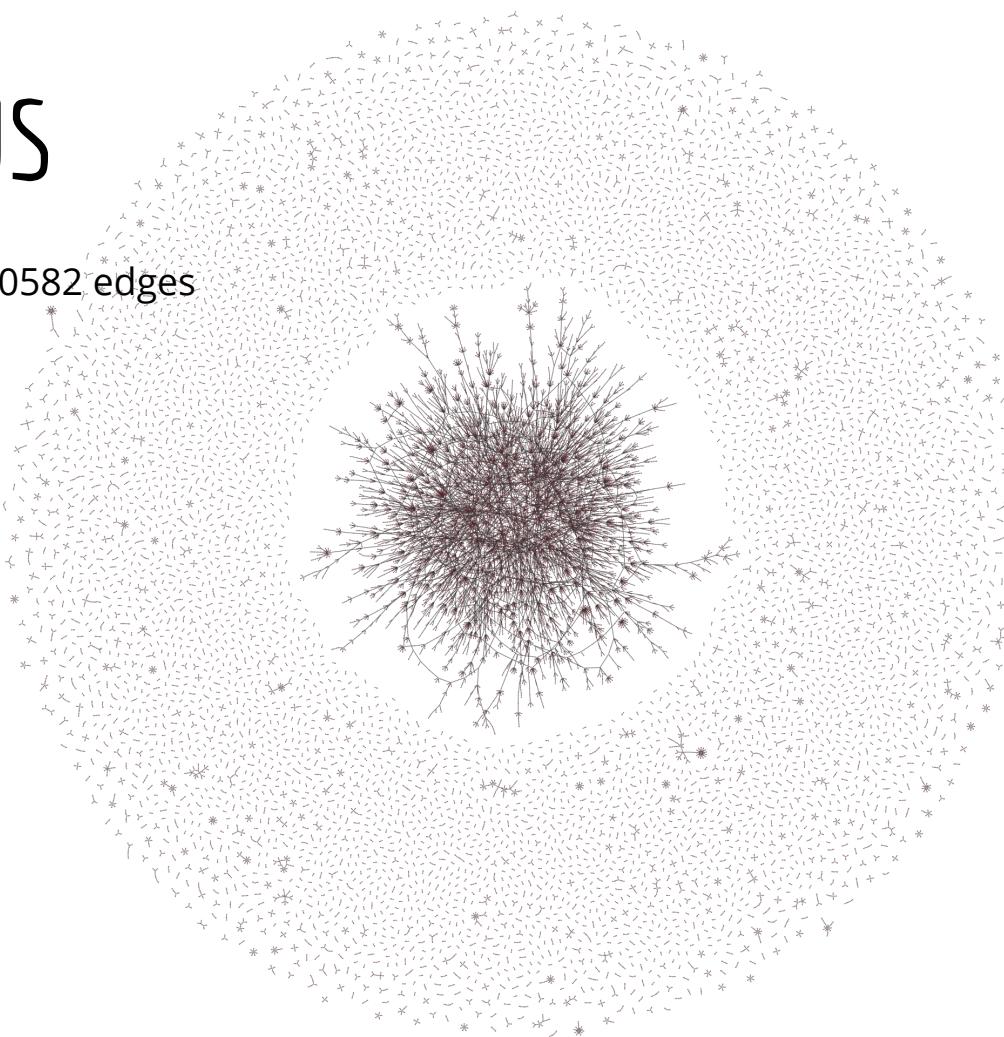
Sequentially selected 23095 edges

(2% of total ratings)



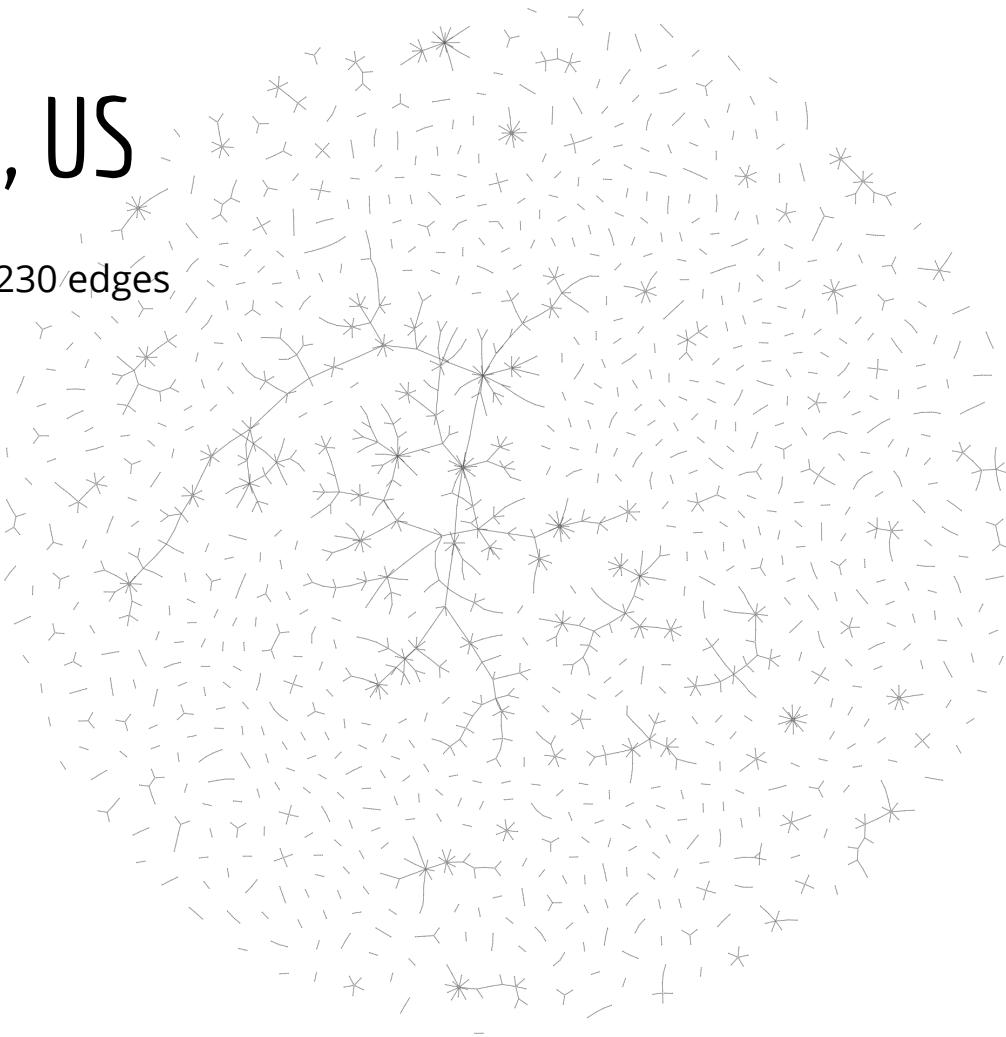
# Phoenix, US

Randomly selected 20582 edges  
(2% of total ratings)

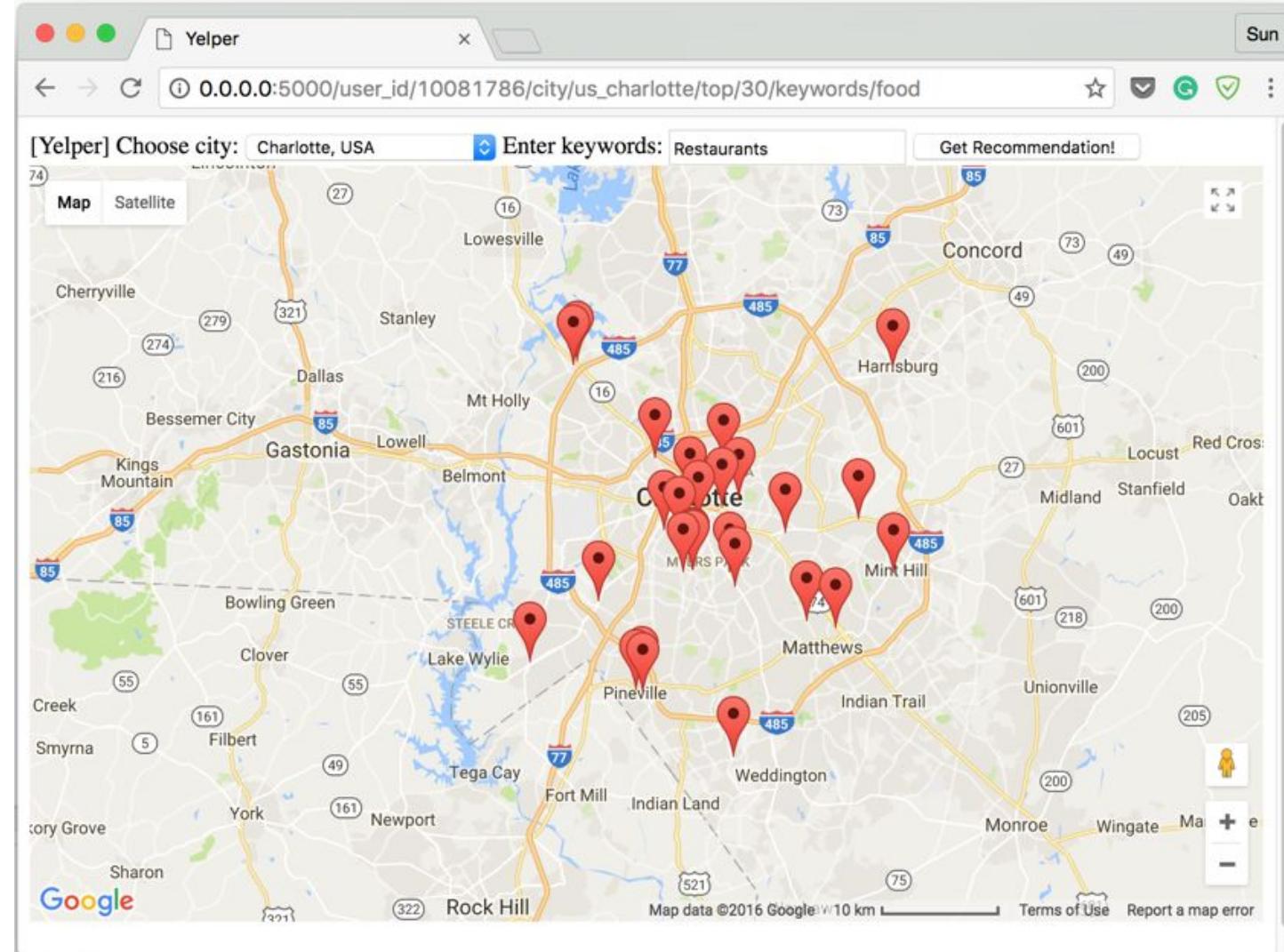


# Pittsburgh, US

Randomly selected 2230 edges  
(2% of total ratings)



Demo 2: Yelper recommendation page



Demo 3: Simulation of user requests handling for recommendation  
using Spark Streaming and Kafka



# Agenda

## What?

What are the components for  
Yelper?

## Why?

Why we need  
recommendation?

## Demo

User-business network; Yelper  
main page; Simulation of  
users' recommendation  
requests handling

## How?

How to build Yelper?

## Summary & Future works

What we learn and what's  
next?

# Summary & future works

- Summary
  - Preprocessing by dividing business data by cities to allow fine tuned and customized recommendations
  - Collaborative filtering based recommendation using Spark MLlib
  - User-business graph visualization using D3 and graph-tool library
  - User-business graph analysis using Spark GraphX in Scala
  - Real-time user request handling simulation using Spark Streaming and Apache Kafka
  - Google Map view to recommend high rated businesses for users
- What we learn
  - Streaming data analysis and mining could be the new norm in the future and we as Data Scientist should be prepared for it
- Future works
  - More graph analysis
    - Graph pagerank analysis using GraphX
    - Community discovery (similar to Facebook social network)
  - Improve recommendation
    - Content-based recommendation
    - Clustering all businesses
    - Extract object from business photos using Convolutional Neural Network

Thanks!