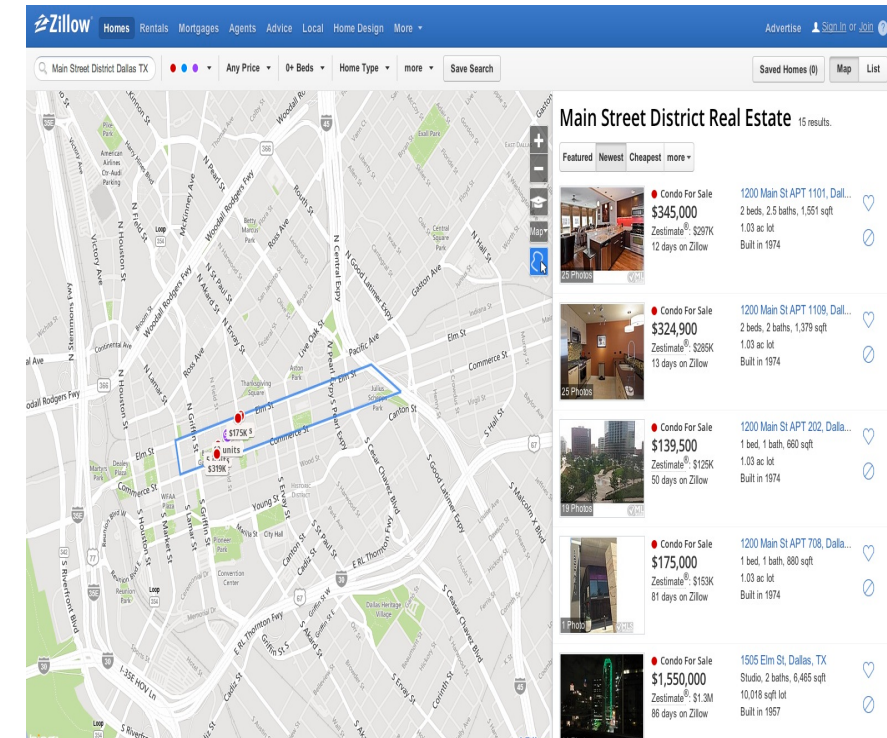# Zillow End-to-End Data Pipeline

BY JERRY KHONG

# Introduction



- Zillow is a real estate marketplace platform

- Buy, sell, rent

- 160 million active monthly users [1]

[1]: HTTP://EXPANDEDRAMBLINGS.COM/INDEX.PHP/ZILLOW-STATISTICS/

galvanizeU
powered by University of New Haven

# Problems in Real Estate

- U.S. residential real estate market is one of the least liquid markets

- Real estate market is volatile and hard to predict

- Real estate transactions are stressful
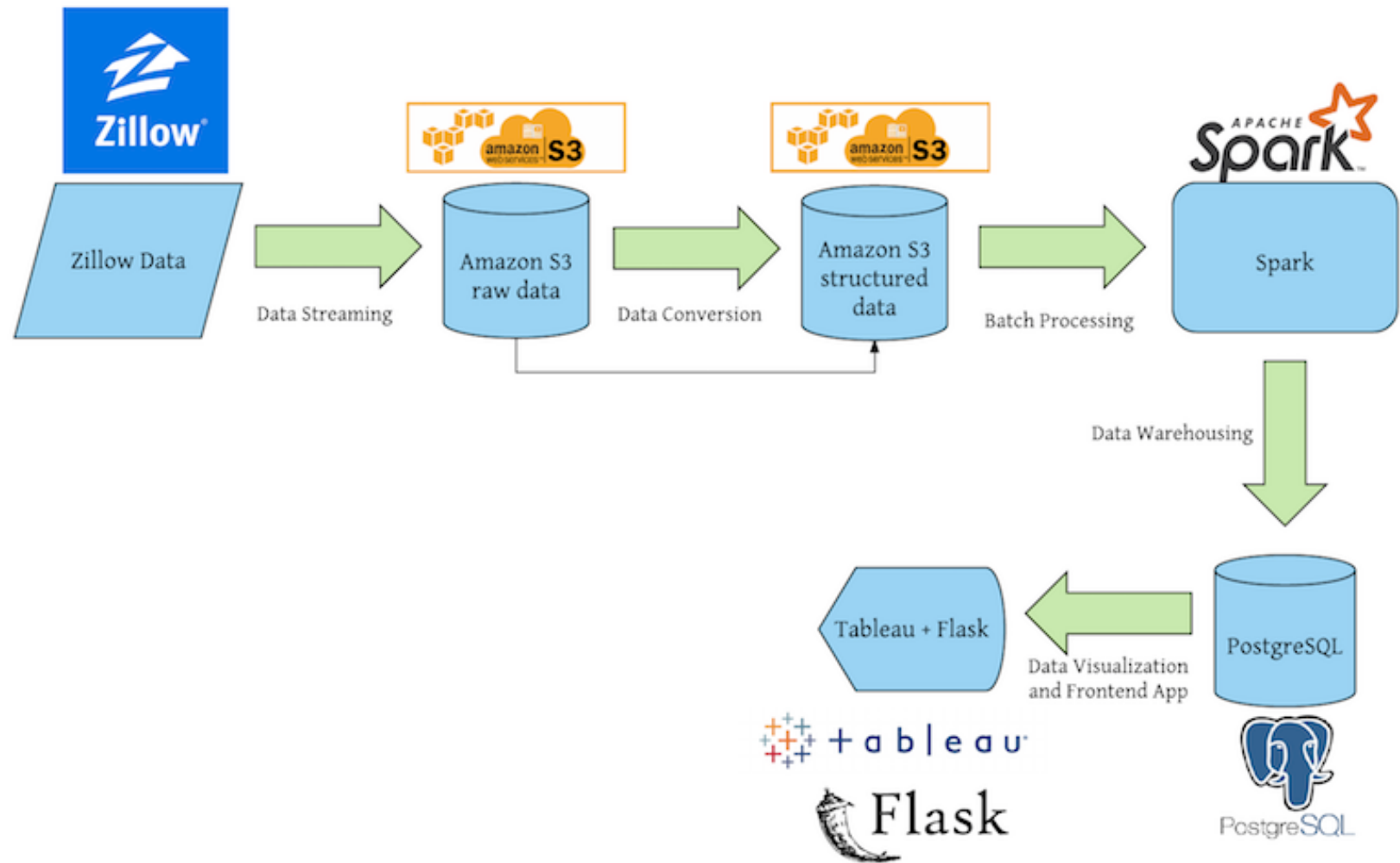
- Zillow's estimated house valuations are bad

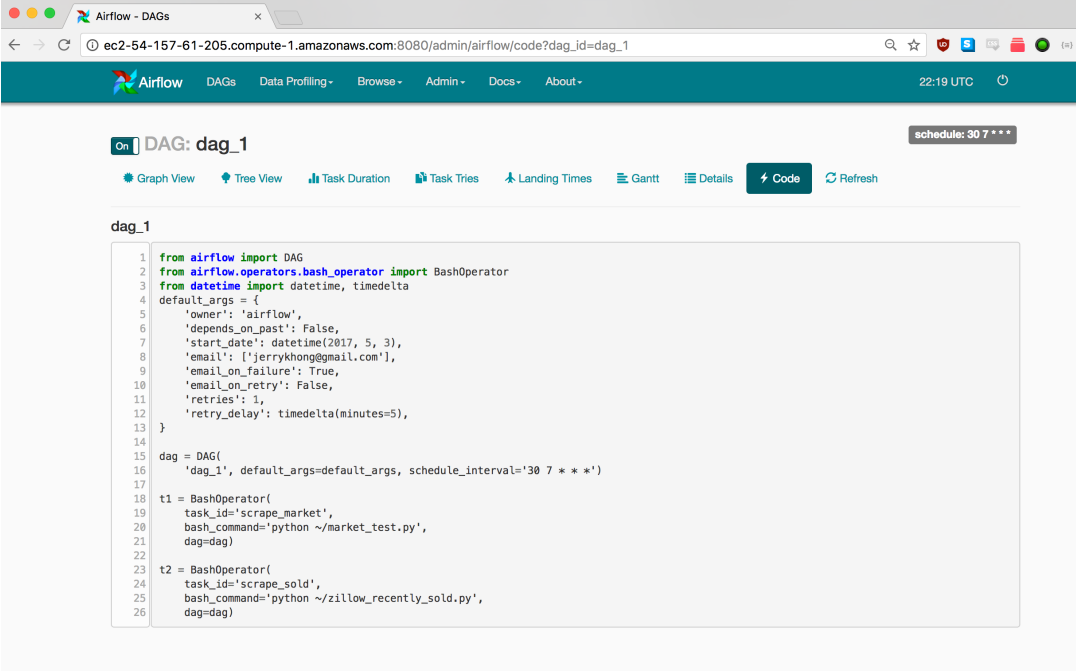# What is the true value of a home?

# Solution

- Assess which areas will be good real estate investments based off historical and current market data.

- Provide users with interactive map to see which zip codes are currently under-valued and over-valued.

- Future product: scale out to display individual listings

# Data Pipeline Architecture

# Stream

- EC2 instance bootstrapped with Selenium and Chrome

- Setup Airflow DAG (daily)
    I.   Scrape all houses on the market
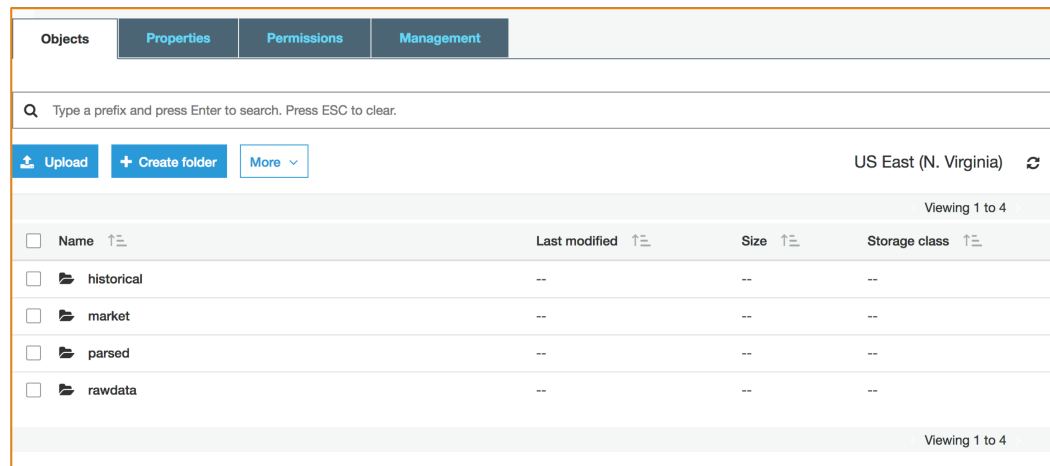    II.  Scrape all recently sold houses

# Store

- Store raw HTML data in S3

- Parse raw HTML data into CSV and store in S3

# Structure

- Use Apache Spark to push parsed data in Postgres DB (3NF)

- Use Apache Spark to push historical data in Postgres DB (3NF)

# Synthesize

- SQL queries for machine learning → export as CSV file

- XGBoost Regression on individual listings to predict price

- Connect Postgres DB to Tableau

- SQL queries in Tableau to filter median house prices

# Show

- Create data visualizations in Tableau

- Embed interactive dashboard on website

# 8 Properties for Big Data systems

**Robustness and Fault Tolerance**

All systems belong to AWS, which is highly robust. In addition, all systems integrates with one another. If the scraping goes down, data itself will be safe in S3.

**Low latency reads and updates**
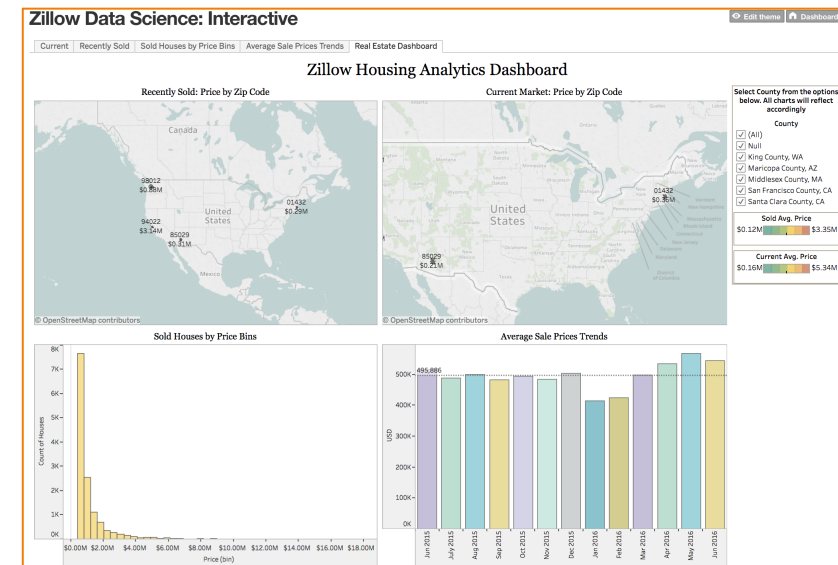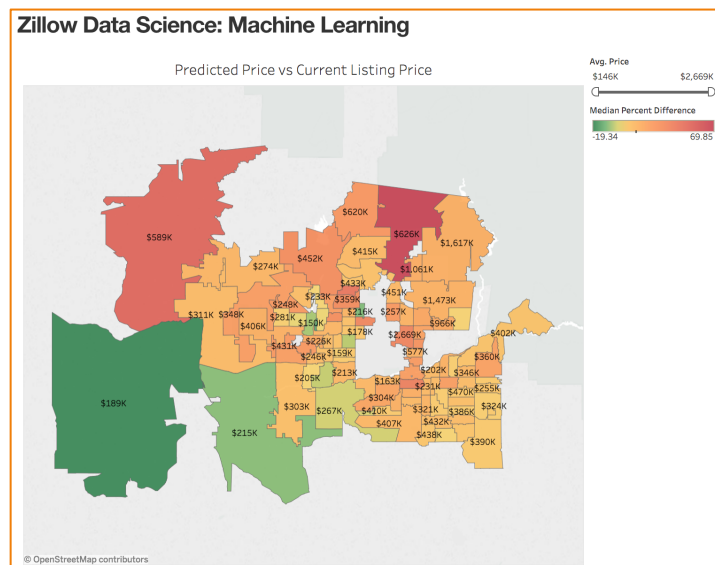
Since the real estate sales cycle is relatively slow, there is no need for real time updates.

**Scalability**

A majority of these technologies are fully scalable. S3 and Spark have high scalability. MySQL or SparkSQL might be a better fit than Postgres for scalability. Tableau is also very scalable.

**Generalization**

This data pipeline architecture is easily extendible for other scraping and can be reused.

# 8 Properties for Big Data systems (cntd)

**Extensibility**

All of the systems in place support the addition of new data sources, features, and models.

**Ad hoc queries**

Postgres and Tableau will be used to do ad hoc queries.

**Minimal maintenance**

This system is complex, and requires monitoring to make sure everything is running. I attempted to make an e-mail airflow DAG in case of errors.

**Debuggability**

Data will always be stored in complete form in S3, so bugs can be easily traced out. Ability for data to be restructured and models can be recomputed if something goes wrong.

# Challenges and Potential Risks

- Working with HTML vs JSON

- Scraping vs API usage

- Bootstrapping different packages

- Getting banned from Zillow

- Research shows that Zillow's data is not up to date (slow)

# Future Goals

- Display individual listings

- Scrape more features

- Incorporate WalkScore and GreatSchools API