**MICHIGAN STATE**
**U N I V E R S I T Y**

# Beta Presentation
## SIFT: *Seller-Forums Information Filtering Tool*

## The Capstone Experience

### Team Amazon

Maxime Goovaerts
Carl Johnson
Luke Pritchett
Benjamin Taylor
Johnny Zheng

Department of Computer Science and Engineering
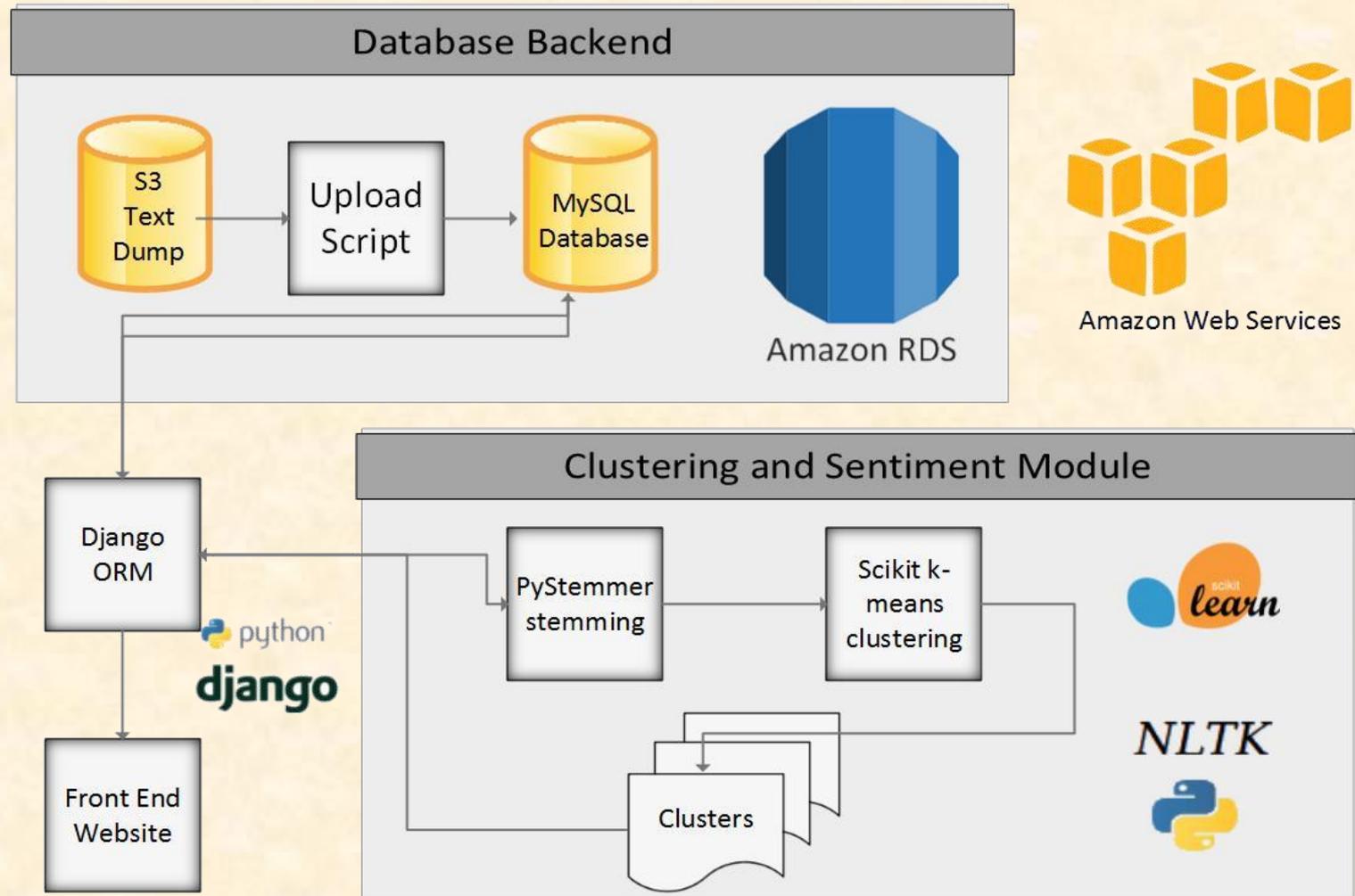Michigan State University

Spring 2015

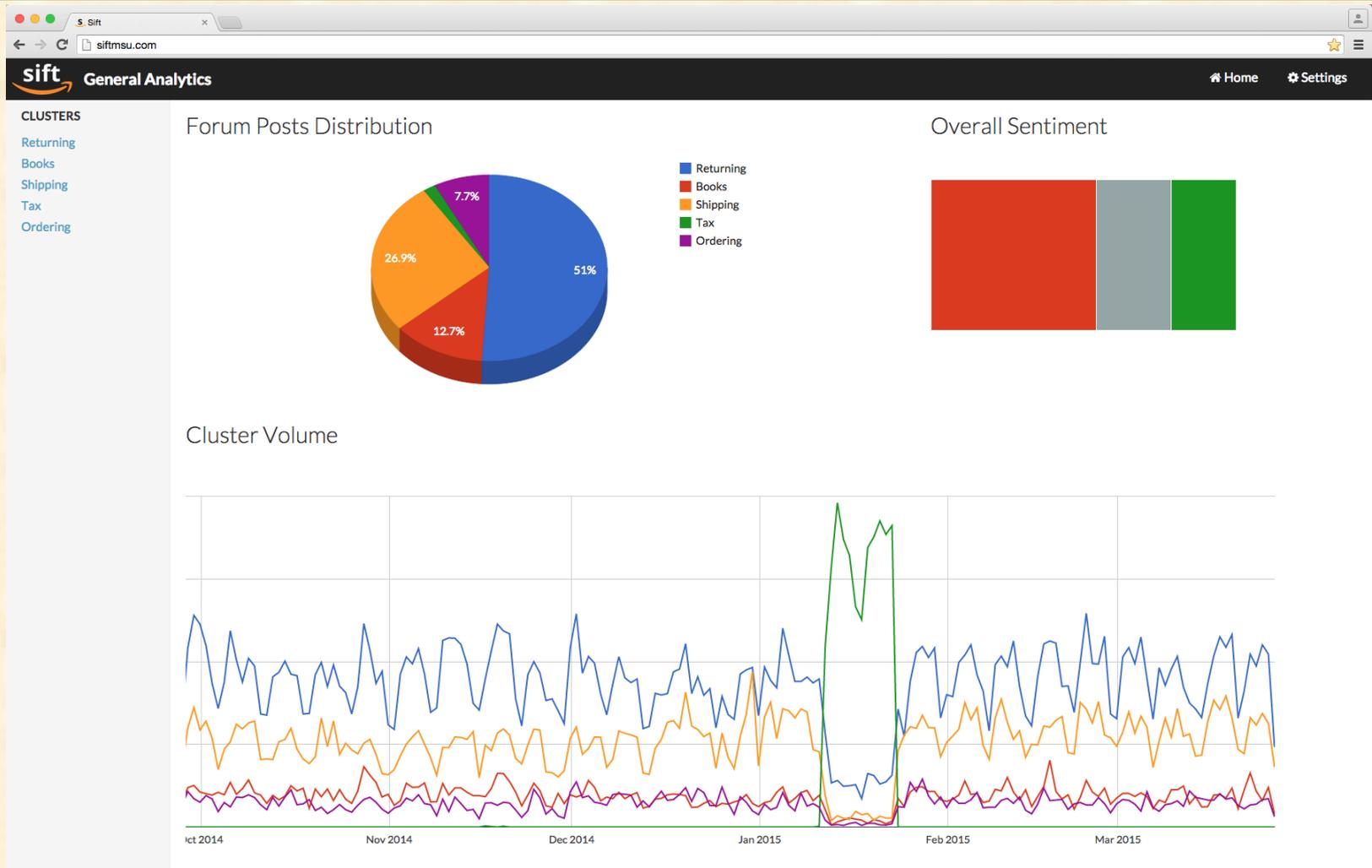*From Students…*
*…to Professionals*

# Project Overview

- Amazon is the largest internet-based retailer
- Unlock the value of 3$^{rd}$ party Seller Forums
- Data Organization and Analysis
  - Clustering
  - Sentiment
- Dashboard
  - Graphs and tables
  - Notifications

# System Architecture

# General Analytics

# Cluster Details

# Settings – *Clusters*

# Settings – *Cluster Configuration*

# Settings – *Notifications*

# What's left to do?

- Diagnostic Clustering Progress Indicator

- Add detail to "scheduled" email

- Testing

- Code Refactoring
  - Speed up page loading by pagination
  - Make sure all functions are properly commented

- Create Project Documentation

- Project Video