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

# Alpha Presentation
## Text Classification of Seller Forums Content

### The Capstone Experience

**Team Amazon**

Maxime Goovaerts
Carl Johnson
Luke Pritchett
Benjamin Taylor
Johnny Zheng

Department of Computer Science and Engineering
Michigan State University
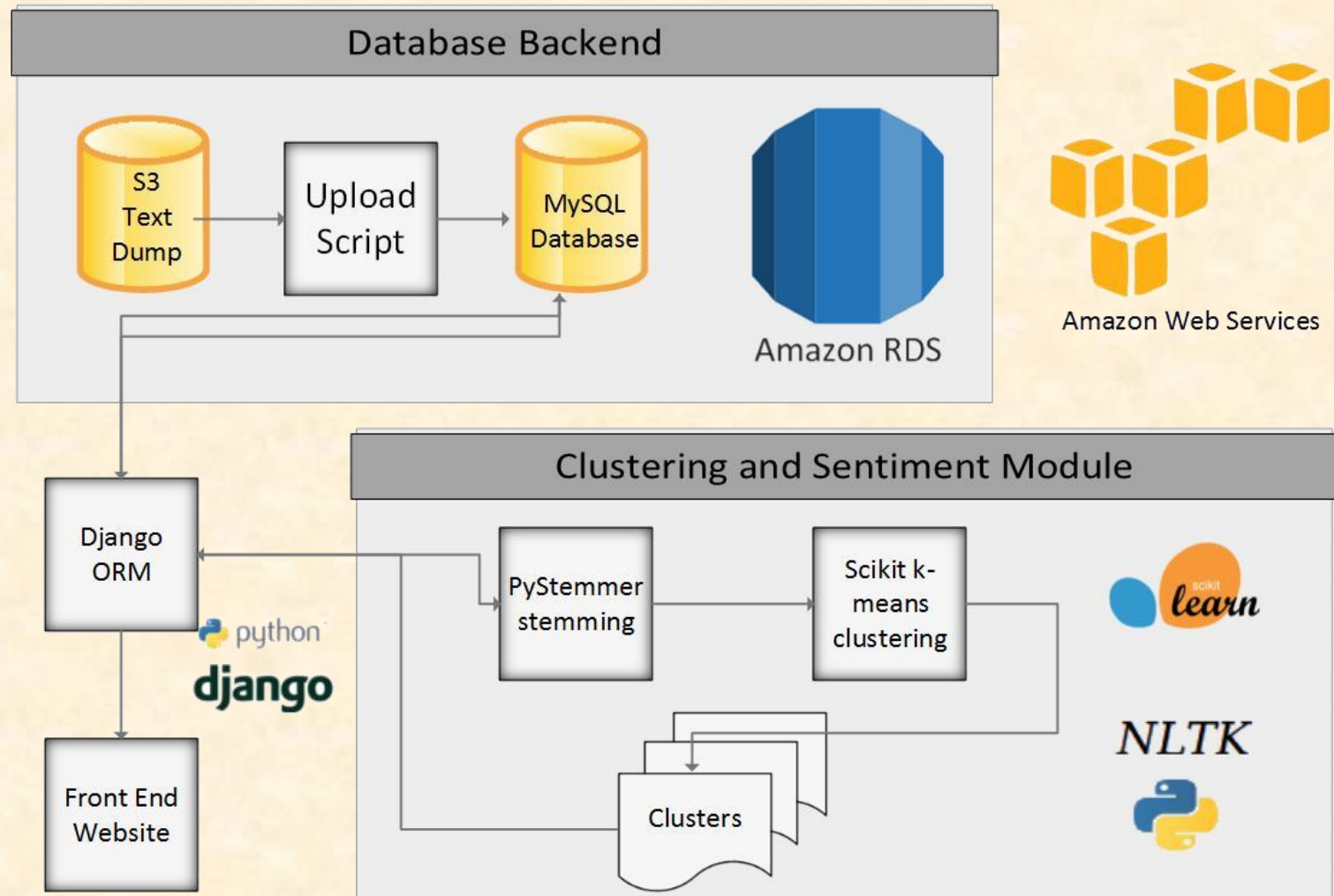
Spring 2015

*From Students…*
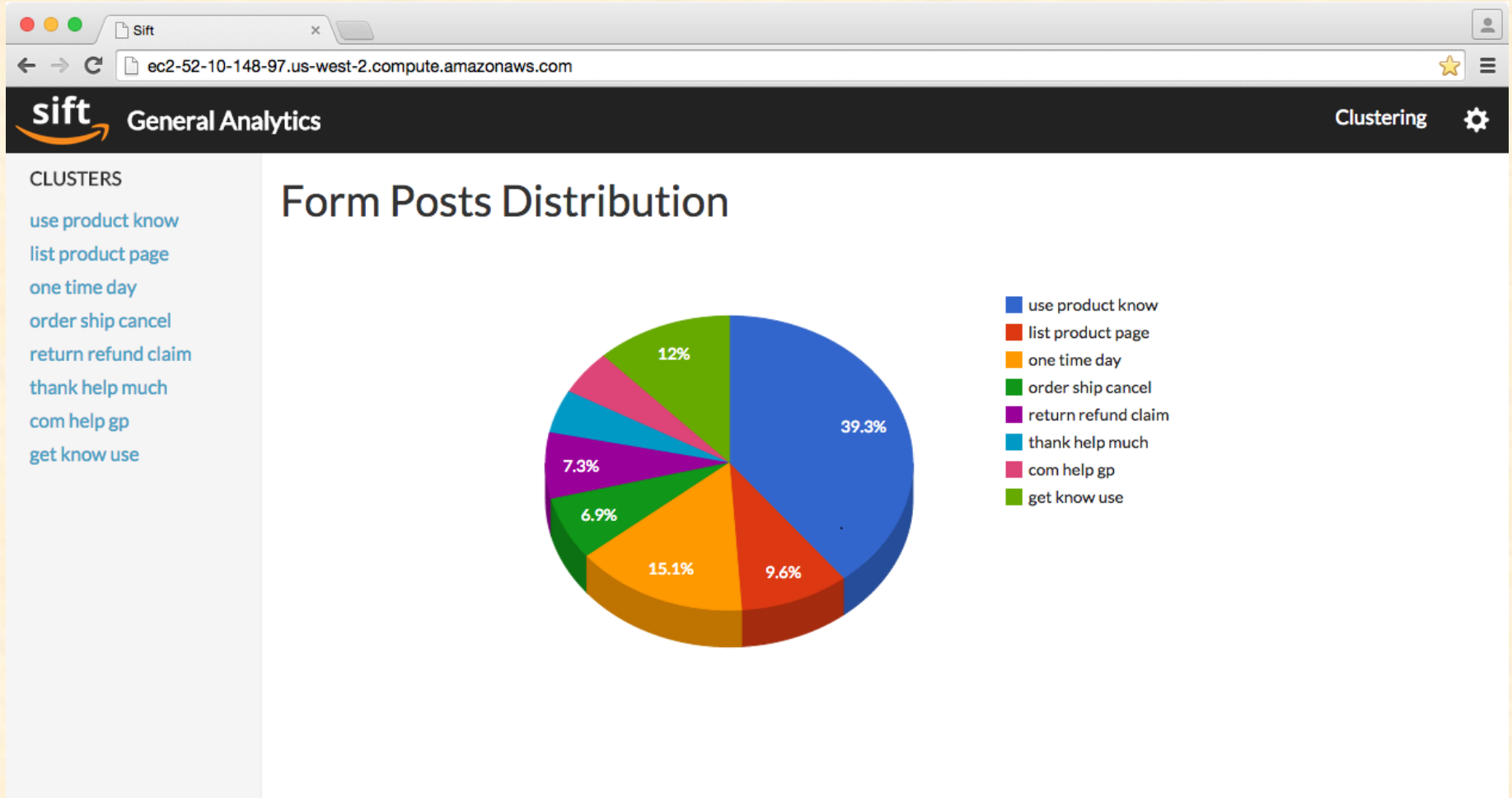*…to Professionals*

# Project Overview

- Amazon is the largest internet-based retailer

- Unlock the value of 3$^{rd}$ party Seller Forums

- Data Organization and Analysis
  - Clustering
  - Sentiment

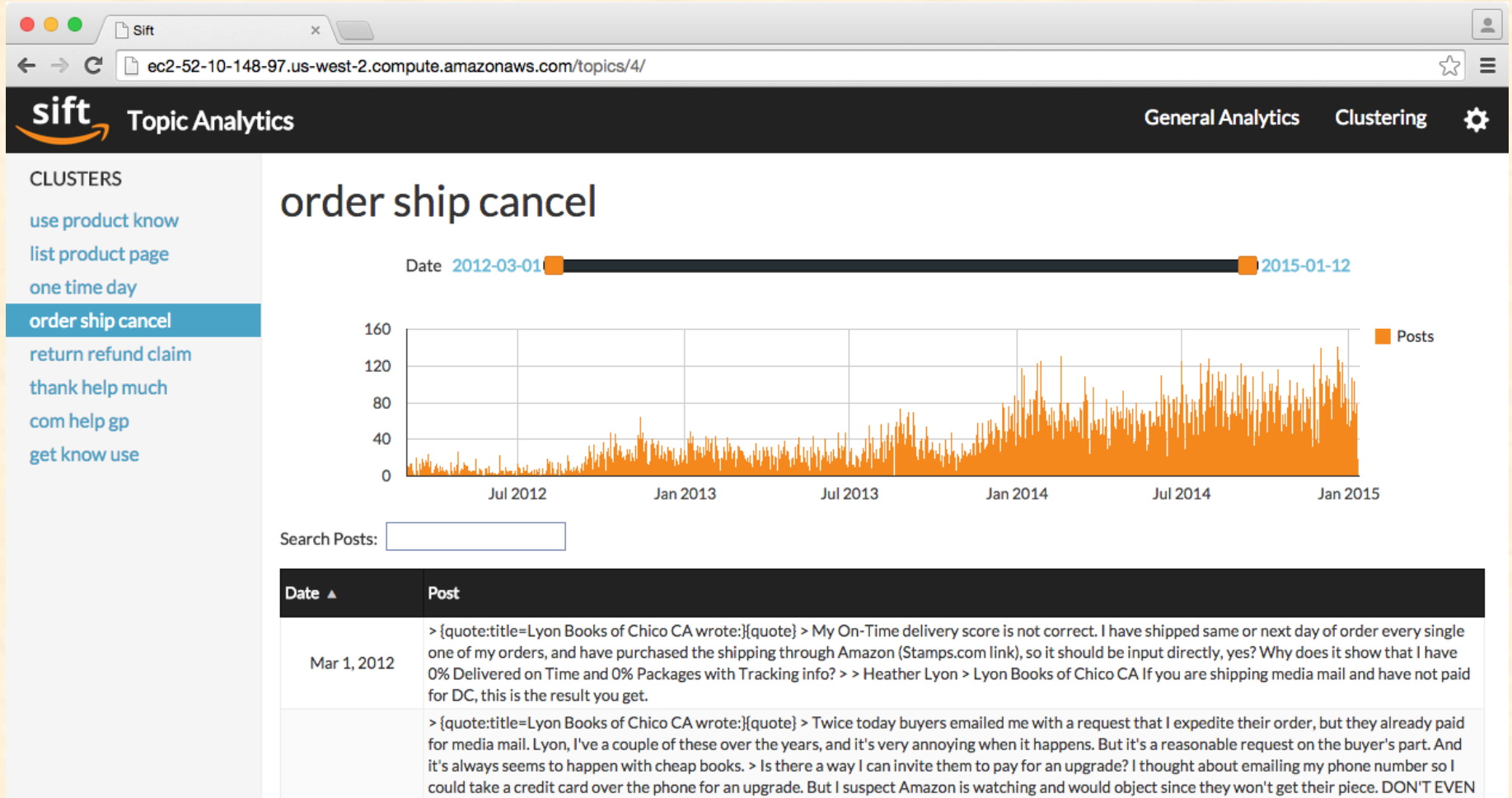- Dashboard
  - Graphs and tables
  - Notifications

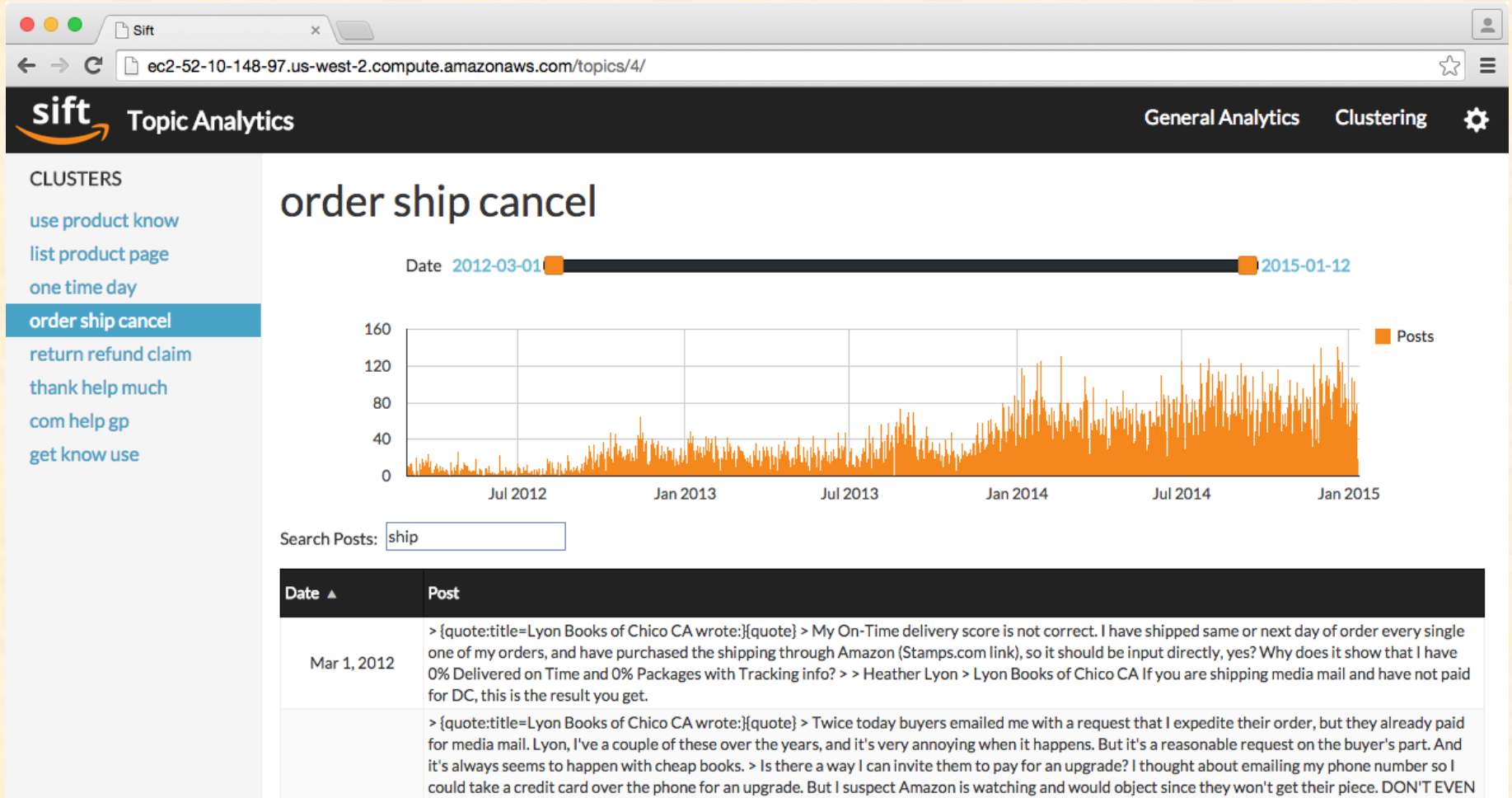# System Architecture

# General Analytics

# Cluster Details

# Cluster Details – *Date Selection*

# Cluster Details – *Searching*

# Clustering Configuration

# What's left to do?

- Optimize Clustering

- Implement Dynamic Loading

- Generate Sci-Kit Cluster Charts

- Track Trending Clusters

- Sentiment Analysis

- Notifications

- Develop more Visualization and Analysis Tools