

CSC:591 - Data Intensive Computing:

Project Status Report

<p style="text-align: center;">GROUP NUMBER 11:</p> <p style="text-align: center;">TweetSmart (Smart Twitter Notification System)</p> <ul style="list-style-type: none"> • TweetSmart is a data-intensive computing system for dynamically determining the top 5 interests of each user and notifying them about the most recent updates regarding those interests. • A data intensive solution is implemented to query this distributed database containing historical and streaming data to calculate those interest values with a few parameters. • This framework will help to analyze topics of interest over time, locations, age of users, etc. and enable smarter notifications and targeting of advertisements. <p>Team members: Saurabh Shanbhag, sshanbh2 Sayali Godbole, ssgodbol Aishwarya Sundararajan, asundar2</p>	<p style="text-align: center;">DELIVERABLES</p> <ol style="list-style-type: none"> Design: <ul style="list-style-type: none"> • Designing the components of system. • Determination of the amount of data (likes * usernames) to be processed and stored. • Determine number of nodes for distributed data storage and retrieval. Getting and Storing Twitter Data: <ul style="list-style-type: none"> • Mining streaming data from Twitter and storing it in the Amazon Dynamo/Cassandra database. Using Hadoop/ZooKeeper for Computation: <ul style="list-style-type: none"> • Distributing the data efficiently and computing top interests (Mapreduce) for each user. Implementation: <ul style="list-style-type: none"> • Getting Top Interests: Retrieving and updating top interests continuously. • Watcher for new posts: determining users to be notified when a person (top interest) posts content. Displaying results/Notification: <ul style="list-style-type: none"> • Notifying user on new interesting post.
<p style="text-align: center;">STATUS</p> <ol style="list-style-type: none"> Design: <ul style="list-style-type: none"> • Finalized Data Structure. • Finalized amount of data and number of nodes (15000 users). Getting and Storing Twitter Data: <ul style="list-style-type: none"> • Using tweepy library, retrieved data about usernames and likes. Using Hadoop/ZooKeeper for Computation: <ul style="list-style-type: none"> • Implemented prototype in python to compute top five interests of a user (to be converted to map and reduce jobs). Implementation: <ul style="list-style-type: none"> • Scaling to multiple users is in progress. • Watcher: <i>pending</i>. Notifying users: <i>pending</i> 	<p style="text-align: center;">ISSUES</p> <ul style="list-style-type: none"> • Getting active usernames for the sample data. • Determining optimal intervals for API requests for new posts and likes. • Access issues for notifications. • Parallelizing/distributing the work among nodes to get faster computation. • Determining the right hosting server from available options. (VCL/AWS) • Validating that the users are notified of the right content by cross-checking the Cassandra database for their topics of interest. • Dealing with the time-lag where new topics of interest are added for a specified user in the database while content about those topics need to be notified to the users simultaneously.