

15-319 / 15-619

Cloud Computing

Recitation 10
March 22nd, 2016

Overview

- **Administrative issues**
Office Hours, Piazza guidelines
- **Last week's reflection**
Project 3.3, OLI Unit 4, Module 15, Quiz 8
- **This week's schedule**
 - 15619 Project – Phase 2- March 30th
 - Quiz 9 – March 25th (Unit 4, Module 16 & 17)
 - Project 3.4 – March 27th

Reminders

- Monitor AWS expenses regularly and tag all resources
 - Check your bill (Cost Explorer > filter by tags).
- Piazza Guidelines
 - Please tag your questions appropriately
 - Search for an existing answer first
- Provide clean, modular and well documented code
 - **Large** penalties for not doing so.
- Utilize Office Hours
 - We are here to help (but not to give solutions)
- Use the team AWS account and tag the 15619Project resources carefully.

Reflection on P3.2 and P3.3

Feedback on code submission



Thread-Safety

While grading P3.2 & P3.3 this week, we have found that many students are not properly implementing thread-safe access to shared data structures.

Review [@2129](#)

OLI Unit 4

Review

OLI Unit 4 : Review

- M14 : Cloud Storage Overview
- M15 : Distributed File Systems (HDFS & Ceph)
- **M16 : NoSQL Database Case Studies**
 - HBase, MongoDB, Cassandra, DynamoDB
- **M17: Cloud Object Storage**

Distributed Databases

- In 2004, Amazon.com began to experience the limits of scale on a traditional web-scale system
- Response was a highly available key-value structured storage system called Dynamo (2007)

Problem	Technique used as solution
Data Sharding	Consistent Hashing
Transient Fault Handling	Sloppy Quorum / Hinted Handoff
Permanent Failure Recovery	Anti-entropy using Merkle trees
Membership and Health Checks	Gossip protocols

- Used in S3, DynamoDB, Cassandra

Distributed Databases

- In 2006, Google published details about their implementation of BigTable
- Designed as a “sparse, distributed multi-dimensional sorted map”
- HBase stores members of “column families” adjacent to each other on the file system - columnar data store

Upcoming Deadlines



- Quiz 9 : Unit 4 - Modules 16, 17

Due : 3/25/2016 11:59 PM Pittsburgh

- Project 3.4 : Social Network with Heterogeneous DBs

Due : 3/27/2016 11:59 PM Pittsburgh

- 15619Project : Phase 2

Due : 3/30/2016 3:59 PM Pittsburgh

Project 3

Review

Project 3 Weekly Modules


- P3.1: Files, SQL and NoSQL
- P3.2: Sharding and Replication
- P3.3: Consistency
- P3.4: Social network with heterogeneous backends
- P3.5: Data warehousing and OLAP

Project 3.4 : Introduction

- Build a social network about movies:

Task 4

Search

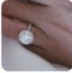



Hacken Lee

81 followers


About View all


Work and Education

 **Diamond Rings** uploaded a photo
2015-09-10 02:47:35



Wow, just experienced Tom and Jerry

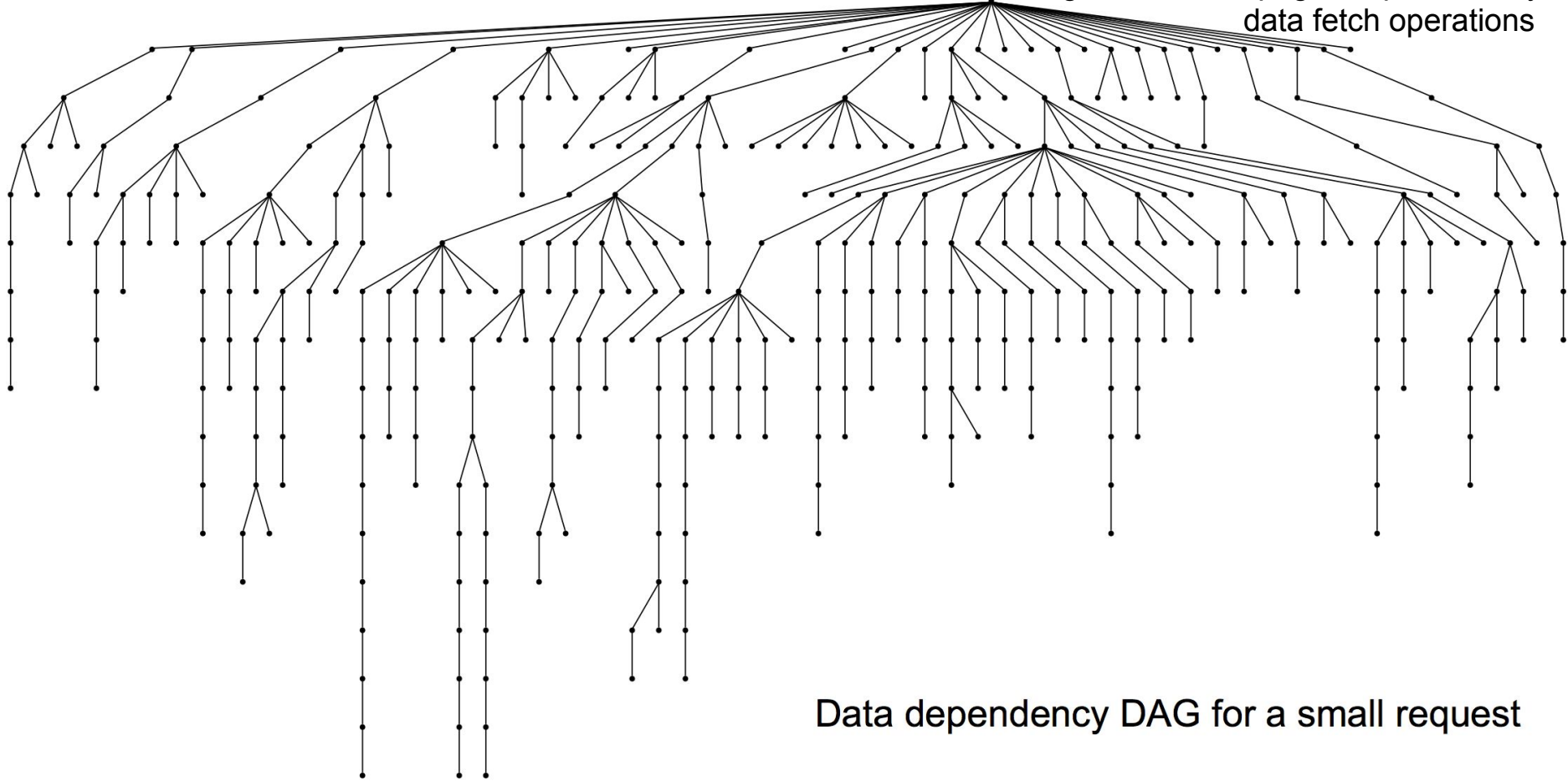
Add a comment 

 **Aimable**
2015-10-04 08:55:40

My mother introduced me to this genre of movies and this was the second one I ever saw. As a non African it takes a while to understand the jargon, but the acting makes up fully for what you do not understand. The two boys in the movie OSITA IHEME; CHINADU IKEDIEZE are so funny and mischievous, their acting alone sold me on getting more featuring them and they haven't disappointed me yet. it is worth seeing I know you will get a laugh out of it.

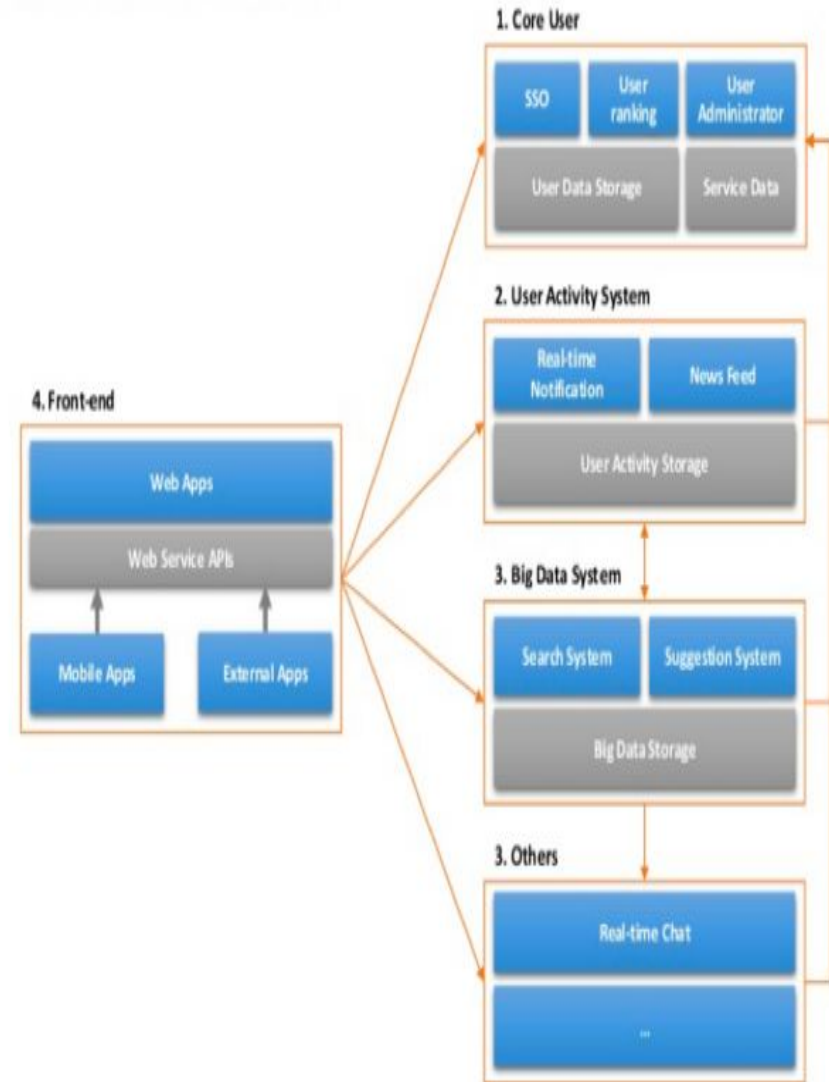
High Fanout and Multiple Rounds of Data Fetching

A single Facebook page, requires many data fetch operations



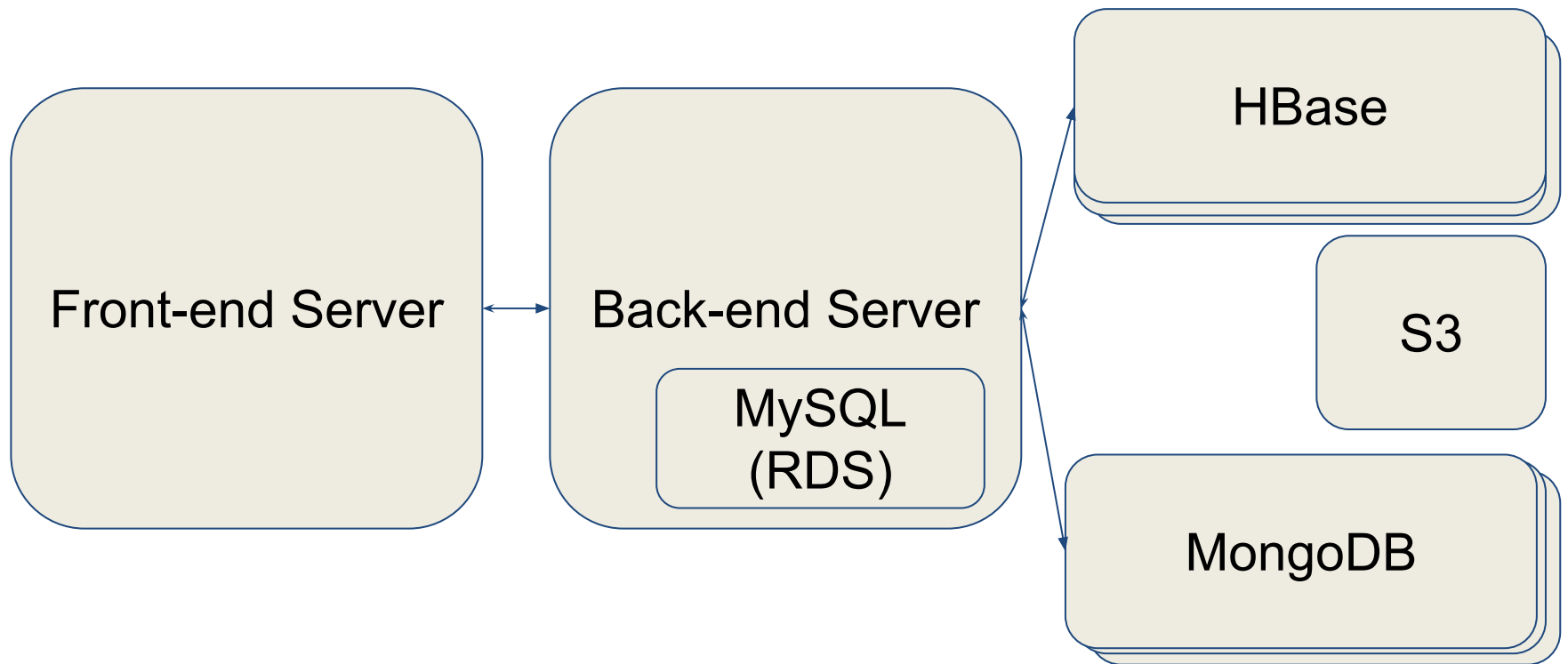
P3.4 Data Set

1. User Profiles
 1. User Authentication System (such as a Single-Sign-On or SSO) - **RDS MySQL**
 2. User Info / Profile - **RDS MySQL**
 3. Action Log
 4. Social Graph of the User: follower, followee, family etc. - **HBase**
2. User Activity System - All user generated media - **MongoDB**
3. Big Data Analytics System
 1. Search System
 2. Recommender System
 3. User Behaviour Analysis



Project 3.4 : Architecture

- **Build a social network about movies:**



MongoDB

- Document Database
 - Schema-less model
- Scalable
 - Automatically shards data among multiple servers
 - Does load-balancing
- Complex Queries
 - MapReduce style filter and aggregations
 - Geo-spatial queries

Heterogeneous Backends

Dataset Name	Database Used	Description	Location
Login Information	MySQL	[UserID, Password]	/home/ubuntu/users.csv
User Profile	MySQL	[UserID, Name, Profile Image URL]	/home/ubuntu/userinfo.csv

Servlet Name	Directory
ProfileServlet.java	/home/ubuntu/Project3_4/src/main/java/cc/cmu/edu/minisite/

Dataset Name	Database Used	Description	Location
Relation	HBase	[Followee, Follower]	/home/ubuntu/links.csv

Servlet Name	Directory
FollowerServlet.java	/home/ubuntu/Project3_4/src/main/java/cc/cmu/edu/minisite/

Servlet Name	Directory
HomepageServlet.java	/home/ubuntu/Project3_4/src/main/java/cc/cmu/edu/minisite/

Dataset Name	Database Used	Description	Location
Posts	MongoDB	Please see the Implementation Requirements section below.	/home/ubuntu/links.csv (MongoDB Instance)

Project 3.4 : Tasks

- **Build a social network about movies:**

- ✓ Task1: Implementing Basic Login with MySQL on RDS

- ✓ Task2: Storing Social Graph using HBase

- ✓ Task3: Build Homepage using MongoDB

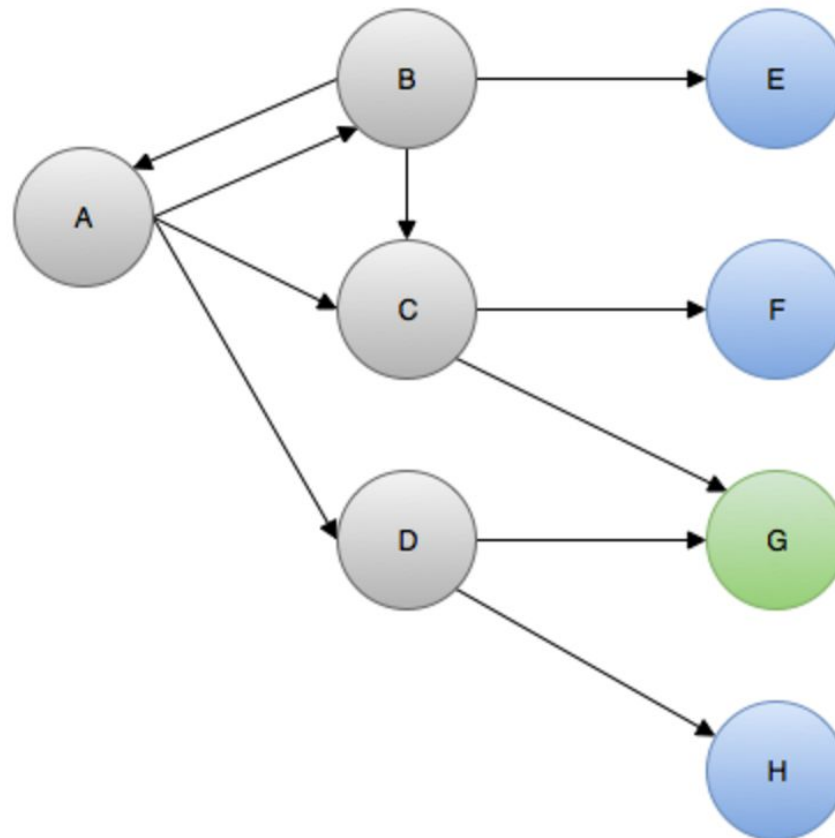
- ✓ Task4: Put Everything Together

- ✓ Bonus Task: Basic Recommendation

Dataset Name	Data Store Used
Login Information	MySQL (RDS)
User Profile	MySQL (RDS)
Relation	HBase
Posts	MongoDB
Profile and Post Images	S3

Project 3.4 : Bonus Task

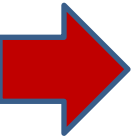
- **Friend recommendation**





Upcoming Deadlines





- Quiz 9 : Unit 4 - Modules 16, 17
Due : 3/25/2016 11:59 PM Pittsburgh
- Project 3.4 : Social Network with Heterogeneous DBs
Due : 3/27/2016 11:59 PM Pittsburgh
- 15619Project : Phase 2
Due : 3/30/2016 3:59 PM Pittsburgh

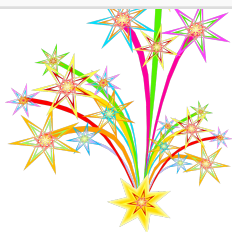


TWITTER DATA ANALYTICS: 15619 PROJECT

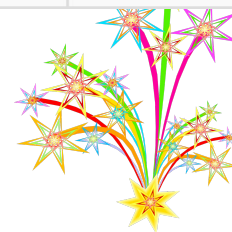


Phase 1 Leaderboard

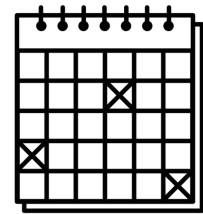
Rank 	Nickname 	Time 	Total 	Effective Throughput 
1	skyleg	03/15/2016 23:28 -0400	200	51752.2
2	Sugoyi	03/14/2016 21:22 -0400	200	51699
3	OnePiece	03/14/2016 23:17 -0400	200	49229.8
4	MyLittlePony	03/15/2016 03:07 -0400	200	47422.36
5	C.C.Lemon	03/16/2016 00:01 -0400	200	47123.4



Well done !!!
Congratulations skyleg & Sugoyi



15619 Project Time Table

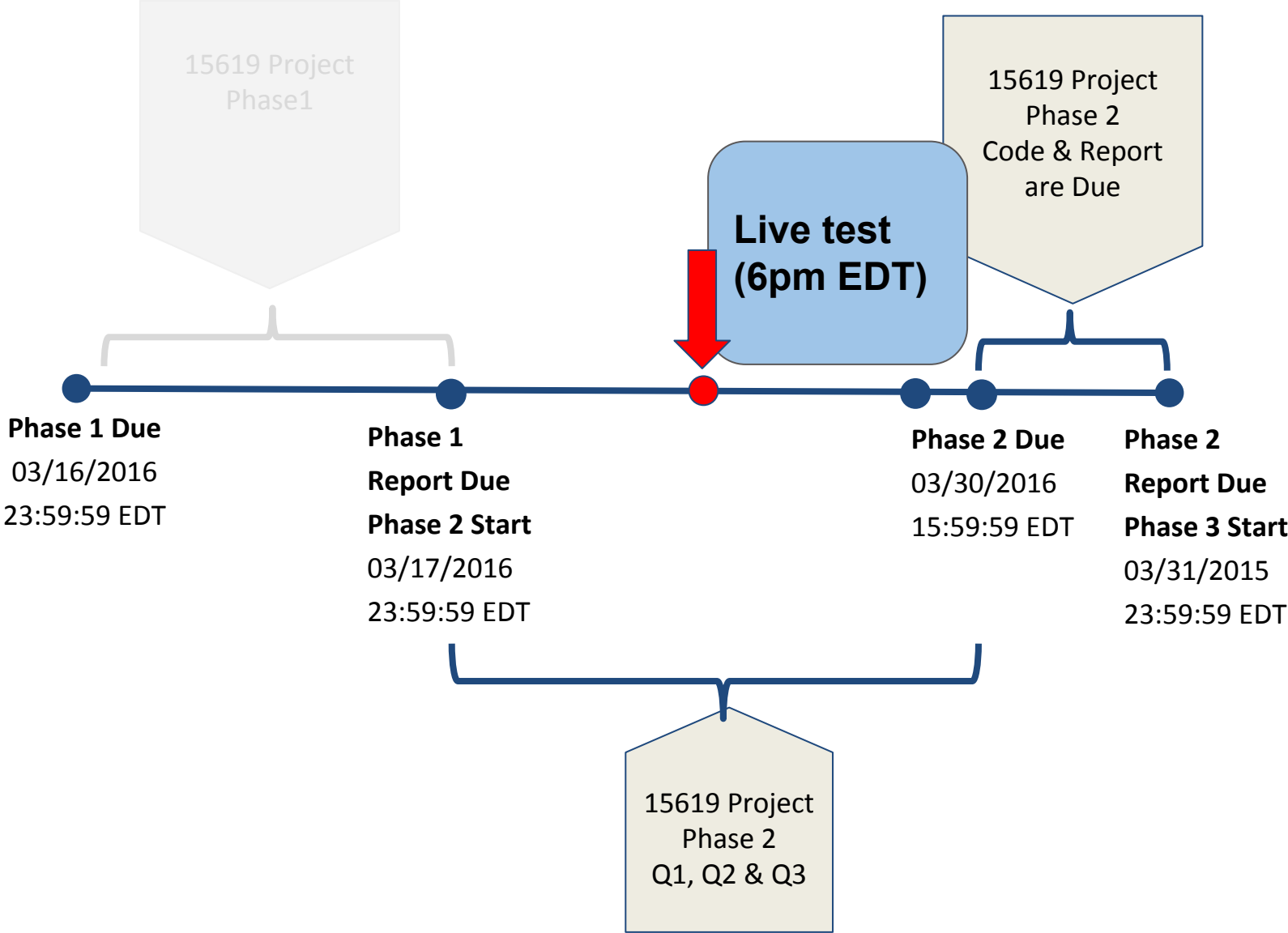


Phase (and query due)	Start	Deadline	Code and Report Due
Phase 1 Part 1 • Q1, Q2	Thursday 02/25/2016 00:00:01 EST	Wednesday 03/16/2016 23:59:59 <u>EDT</u>	Thursday 03/17/2016 23:59:59 <u>EDT</u>
Phase 2 • Q1, Q2, Q3	Thursday 03/17/2016 00:00:01 <u>EDT</u>	Wednesday 03/30/2016 15:59:59 <u>EDT</u>	
Phase 2 Live Test (Hbase/MySQL) • Q1, Q2, Q3	Wednesday 03/30/2016 18:00:01 <u>EDT</u>	Wednesday 03/30/2016 23:59:59 <u>EDT</u>	Thursday 03/31/2016 23:59:59 <u>EDT</u>
Phase 3 • Q1, Q2, Q3, Q4	Thursday 03/31/2016 00:00:01 <u>EDT</u>	Wednesday 04/13/2016 15:59:59 <u>EDT</u>	
Phase 3 Live Test • Q1, Q2, Q3, Q4	Wednesday 04/13/2016 18:00:01 <u>EDT</u>	Wednesday 04/13/2016 23:59:59 <u>EDT</u>	Thursday 04/13/2016 23:59:59 <u>EDT</u>

Note:

- There will be a report due at the end of each phase, where you are expected to discuss design, exploration and optimizations.
-

15619 Project Phase 2 Deadlines



Phase 2

- One new query (Q1, Q2 and **Q3**)
 - More ETL
 - Multiple tables and queries

- Live Test!!!
 - Both HBase and MySQL
 - Two DNS
 - Includes Mixed-Load

Hints - MySQL

- System Environment
 - Storage Medium
 - Storage Engine
 - Character set
 - Import data (SHOW WARNINGS)
 - Indexing
- Profiling/Optimization
 - EXPLAIN
 - SET PROFILING=1
 - htop, iotop

Hints - HBase

- Loading data:
 - Pig, thrift, MapReduce
- HBase schema:
 - GET is much faster than SCAN
 - How to design rowkey?
- HBase cluster:
 - Cloudera Manager - easy deployment and management of cluster
 - Deploy your own HBase cluster and automate it
 - Using EMR will lead to higher cost ⇒ must use less instances <\$.85
- HBase configuration tuning:
 - Heap size, cache size
 - Block size
 - Region size / number
 - The book: <http://shop.oreilly.com/product/0636920014348.do>

<http://archive.cloudera.com/cdh5/cdh/5/hbase-0.98.1-cdh5.1.5/book/ops.capacity.html>

General Tips

- Don't blindly optimize for every component, identify the bottlenecks using fine-grained profiling
- Use caches wisely: cache in HBase and MySQL is obviously important, but front-end cache will most likely to fail during the Live test
- Get the whole picture of the database you are using, don't just Google and adopt "HBase/MySQL optimization techniques" blindly.
- Review what we have learned in previous project modules
 - Scale out
 - Load balancing
 - Replication and sharding
- Look at the feedback of your Phase 1 report!

Q3 : Handling Complex Read Queries

- Calculate word occurrences in tweet text within a certain user id range and a date range. (Two-range query)

- Request Format

```
GET/q3?start_date=yyyy-mm-dd&end_date=yyyy-mm-dd&start_userid=uid&end_userid=uid&words=w1,w2,w3
```

- Response Format

```
TEAMID,TEAM_AWS_ACCOUNT_ID\n
```

```
w1:count1\n
```

```
w2:count2\n
```

```
w3:count3\n
```

- Target RPS **6000**

Q3 : Handling Complex Read Queries

- Request Example (Double Range Query)

```
GET/q3?start_date=2014-04-01&end_date=2014-05-28&start_userid=51538630&end_userid=51539182&words=u,petition,loving
```

- Response Format

```
Team,1234-5678-1234
```

```
u:7\n
```

```
petition:2\n
```

```
loving:5\n
```

Q3: ETL

1. Filter out non-english tweets (lang != 'en')
2. Split words when a non-alphanumeric character ([^a-zA-Z0-9]) is encountered. Use the regular expression provided in the write-up.
2. Words are case **INSENSITIVE** in word count.
3. Banned words in Q2 will not appear in Q3 requests.
4. Ignore words from stop words list.
5. Match your ETL result with the reference file provided.

We will provide a Q3 reference file and reference server.

Q3 Hints

- Evaluate the Q3 functionality when designing your schema, especially for HBase. Your schema design should avoid a big scan.
- Try to get an idea of the size of user id ranges and date ranges in the requests.

Phase 2 Live Test

HBase LiveTest

Time	Value	Target	Weight
6:00 pm - 6:30 pm	Warm-up (Q1 only)	-	0%
6:30 pm - 7:00 pm	Q1	27000	5%
7:00 pm - 7:30 pm	Q2	10000	10%
7:30 pm - 8:00 pm	Q3	6000	10%
8:00 pm - 8:30 pm	Mixed Reads(Q1,Q2, Q3)	TBD	5+5+5 = 15%

Half Hour Break

MySQL LiveTest

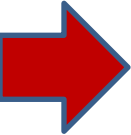
Time	Value	Target	Weight
9:00 pm : 9:30 pm	Warm-up	-	0%
9:30 pm - 10:00 pm	Q1	27000	5%
10:00 pm - 10:30 pm	Q2	10000	10%
10:30 pm - 11:00 pm	Q3	6000	10%
11:00 pm - 11:30 pm	Mixed Reads (Q1,Q2,Q3)	TBD/TBD/TBD	5+5+5 = 15%

Tips for Phase 2

- Carefully design and check ETL process
- Watch your budget:
 - \$60 = Phase 2 + Live Test
- Preparing for the live test
 - You are required to submit two URLs, one for the MySQL live test and one for the HBase for live test
 - Budget limited to \$.85/hr for MySQL and HBase web service separately.
 - Need to have all Q1-Q3 running at the same time.
 - Don't expect testing in sequence.
 - Queries will be mixed.



Upcoming Deadlines

- Quiz 9 : Unit 4 - Modules 16, 17
Due: 3/25/2016 11:59 PM Pittsburgh
- Project 3.4 : Social Network with Heterogeneous DBs
Due: 3/27/2016 11:59 PM Pittsburgh
- **15619Project : Phase 2**
 **Due: 3/30/2016 3:59 PM Pittsburgh**