

15-319 / 15-619

Cloud Computing

Recitation 6

February 21th, 2017

Overview

- **Announcements**
- **Last week's reflection**
- Project 2.2, OLI unit 3 module 7, 8 and 9
- **This week's schedule**
 - Unit 3 - Modules 10, 11 and 12
 - Quiz 5 - Friday, February 24th
 - Project 3.1 - Sunday, February 26th

Announcements

- Use spot instances as much as possible
 - e.g. P3.1 HBase cluster
- **Protect your credentials**
 - Crawlers are looking for AWS credentials on public git repos!
- Get started early!!!

Last Week's Reflection

- OLI: Conceptual Content
 - Unit 3 - Modules 7, 8 and 9:
 - Introduction and Motivation
 - Virtualization
 - Resource Virtualization - CPU
 - Quiz 4 completed
- P2.2: Containers: Docker + Kubernetes

Project 2.2

Containers: Docker + Kubernetes

- Build a multi-cloud service to compile and run user code submitted through a front end.
- Project Tasks:
 - Task 1: Simple Docker example.
 - Task 2: Clustering containers in GCP with Kubernetes.
 - Task 3: Code Execution as a Service
 - Task 4: Container Orchestration of xCloud Clusters

Project 2.2 Feedback

[See this Piazza post to give feedback!](#)

Please leave us feedback!!!

This Week: OLI Content

- UNIT 3: Virtualizing Resources for the Cloud
 - Module 7: Introduction and Motivation
 - Module 8: Virtualization
 - Module 9: Resource Virtualization - CPU
 - **Module 10: Resource Virtualization - Memory**
 - **Module 11: Resource Virtualization – I/O**
 - **Module 12: Case Study**
 - Module 13: Storage and Network Virtualization

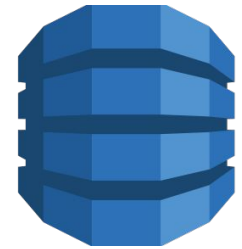
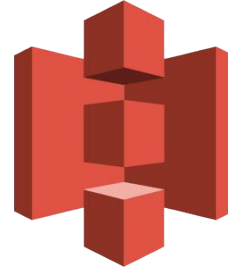
OLI, Unit 3: Modules 10, 11, 12

- Understand two-level page mappings from virtual to real to physical
- Learn how memory is overcommitted and reclaimed using ballooning
- Study how I/O requests are intercepted at different interfaces
- Map these concepts into your practical exploration with AWS

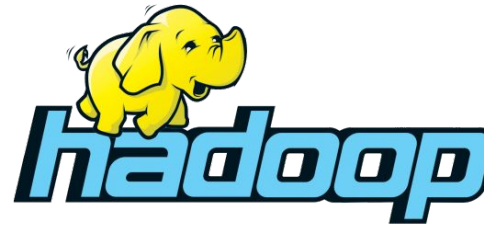
This Week's Project

- P3.1: Files and Databases
 - Comparison and Usage of Flat files, MySQL and HBase
- P3.2: Social Networking Timeline with Heterogeneous Backends
 - Social Networking Timeline with Heterogeneous Backends (MySQL, HBase, MongoDB, S3)
- P3.3: Replication and Consistency
 - Multi-threaded Programming and Consistency

Project 3 - Storage

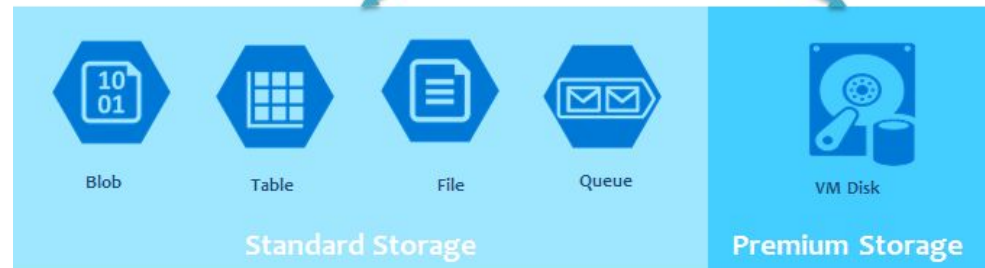


APACHE
HBASE



Azure Storage Account

PostgreSQL



Project 3 Weekly Modules

- P3.1: Files, SQL and NoSQL
 - Storage & IO Benchmarking
 - NoSQL Primer
- P3.2: Social network with heterogeneous backend storage
 - NoSQL Primer
- P3.3: Replication and Consistency models
 - Primer: Intro. to Java Multithreading
 - Primer: Thread-safe programming
 - Primer: Intro. to Consistency Models

Performance Benchmarks

- Running sysbench and preparing data
 - Remember to change to mounted directory.
 - Use the prepare option to generate the data.
- Experiments
 - Run sysbench with different storage systems and instance types.
 - Doing this multiple times to reveal different behaviors and results.
- Compare the requests per second.

Performance Benchmarks Sample Report

Scenario	Instance Type	Storage Type	RPS Range	RPS Increase Across 3 Iterations
1	t1.micro	EBS Magnetic Storage	100, 100.5, 100	Trivial (<5%)
2	t1.micro	EBS General Purpose SSD	617.98, 643.99, 634.33	Trivial (<5%)
3	m3.large	EBS Magnetic Storage	309.88, 383.47, 460.89	Significant (can reach 50% with absolute increase of 150-200)
4	m3.large	EBS General Purpose SSD	1367.32, 1653.04, 1722.52	Noticeable (can reach 25% with absolute increase of 300-400)

What can you conclude from these results?

Performance Benchmarks Conclusion

- SSD has better performance than magnetic disk
- m3.large instance has better performance than t1.micro instances
- The RPS increase across 3 iterations for m3.large is more significant than that for t1.micro:
 - The reason is an instance with more memory can cache more of the previous requests for repeated tests.
 - **Caching is also a vital performance tuning mechanism when building high performance applications.**

NoSQL

- Non-relational or Non-SQL
- Why NoSQL?
- Exercises in the primer
- Motivations for NoSQL
 - ACID vs BASE ...
 - **A**tomic **C**onsistent **I**solated **D**urable
 - **B**asic **A**vailability **S**oft-state **E**ventually Consistent
- NewSQL?

Project 3.1 Overview

P3.1: Files, SQL, and NoSQL:

- Run basic Unix commands like grep, awk etc to extract certain data from given datasets in flat files
- Use relational databases (MySQL)
 - load data, run basic queries
- Use a NoSQL database (HBase)
 - load data, run basic queries

The NoSQL Primer is vital to completing P3.1.

Flat Files

- Flat files, plain text or binary
 - comma-separated values (CSV) format: Carnegie,Cloud
Computing,A,2017
 - tab-separated values (TSV) format: Carnegie\tCloud
Computing\tA\t2017
 - a custom and verbose format: Name: Carnegie, Course:
Cloud Computing, Section: A, Year: 2017
- Lightweight
- Flexible
- Performing analysis on data in files can be expensive
- Insert is expensive due to having to seek
- Managing a large number of files with different data types can become complex ...

Databases

- A collection of organized data
- Database management system (DBMS)
 - Interface between user and data
 - Store/manage/analyze data
- Relational databases
 - Based on the relational model (schema):
MySQL
- NoSQL Databases
 - Unstructured/semi-structured
 - DynamoDB, HBase, Mongo,
Google BigTable

Databases

- Advantages

- Logical and physical data independence
- Concurrency control and transaction support
- Easy and convenient querying (e.g. SQL)
- ...

- Disadvantages

- Cost (computational resources, fixed schema)
- Maintenance and management
- Complex and time-consuming to design schema
- ...

Files vs. Databases

- Compare flat files to databases
- Think about:
 - What are the advantages and disadvantages of using flat files or databases?
 - In what situations would you use a flat file or a database?
 - How to design your own database? How to manipulate and query data in a database?

Flat File Tasks

- Analyze Yelp's Academic Dataset
 - https://www.yelp.com/dataset_challenge
- Answer questions in runner.sh
 - Use tools such as awk, grep, ...
 - Similar to what you did in Project 1.1, 1.2
- Merge csv files by joining on a common field
 - Identify the disadvantages of flat files

MySQL Tasks

- Prepare tables
 - The script to create the table is already provided
 - Find a way to load the data properly into MySQL
- Use MySQL queries to answer questions
 - Learn JDBC
 - Complete MySQLTasks.java
 - Aggregate functions, joins
 - Statement and PreparedStatement
 - SQL injection
- Learn how to use indexes to improve performance

HBase (NoSQL) Tasks

- Launch an EMR cluster with HBase installed.
- Follow the write-up to download and load the data into HBase.
- Try different querying commands in the HBase shell.
- Complete HBaseTasks.java using HBase Java APIs.
- Use Sqoop

Project 3.1 Manual Grading Preview




- To evaluate how well other people can read your code, we will be manually grading your submitted code
 - To enhance readability
 - Use the [Google Code Style](#) guidelines
 - Always add comments especially for complex parts
 - Checkstyle tool is also introduced for you to use

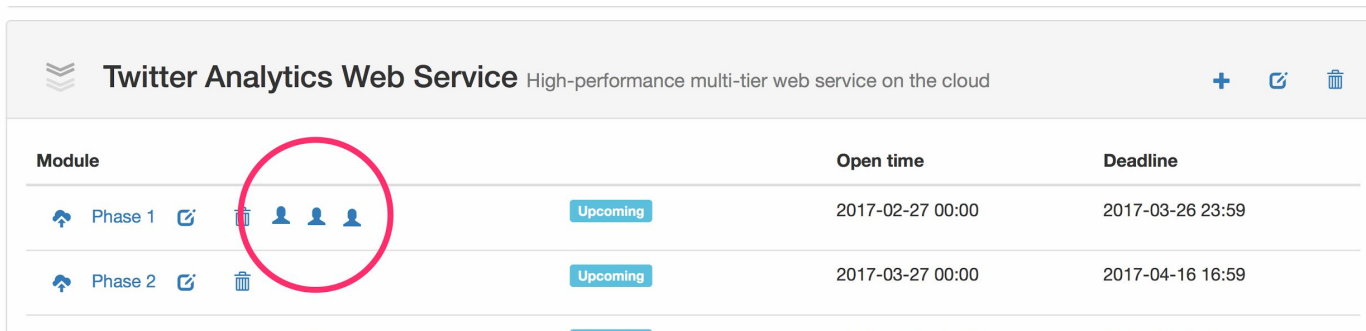
P3.1 Reminders

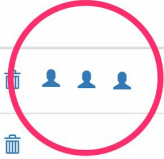
- Tag your resources with: **Key: Project, Value: 3.1**
- You can save a copy of your runner.sh if you want to take a break and work on it later.
- Make sure to terminate the instance after finishing all questions and submitting your answers.

Upcoming Deadlines

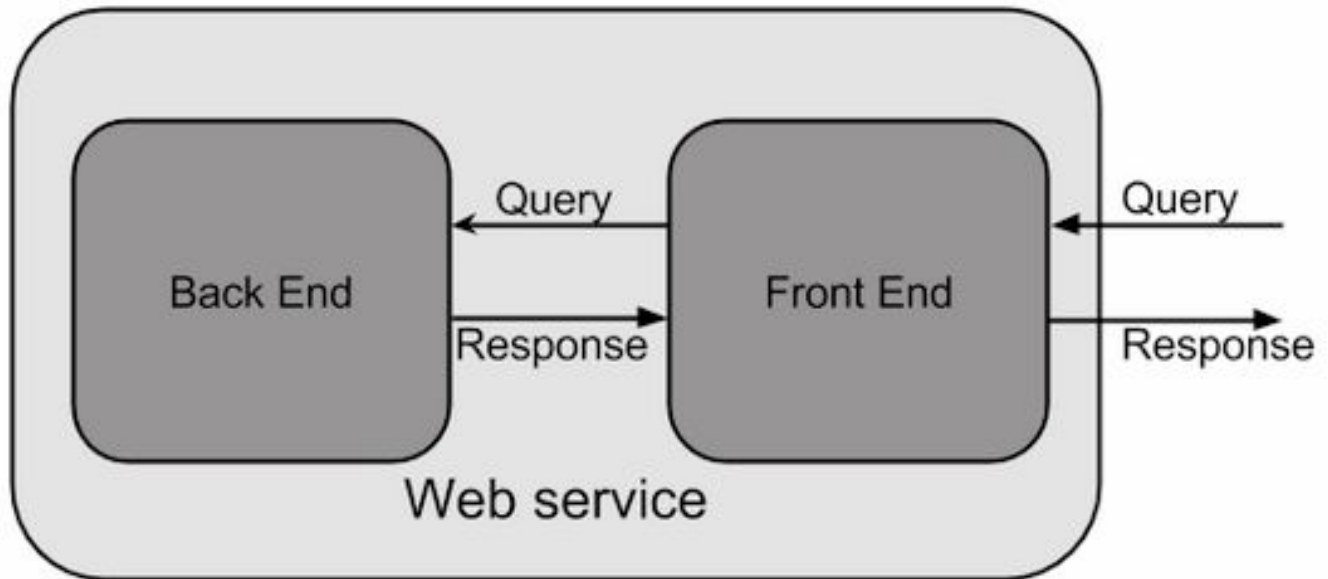


- Quiz 5: Unit 3 - Virtualizing IO, Memory; Cases
Due: **Feb 25 2017 11:59PM Pittsburgh** 
- Project 3.1: Files and Databases
Due: **Feb 26 2017 11:59PM Pittsburgh** 
- Team Project - Team Formation on theproject.zone
Due: **Feb 20 2017 11:59PM Pittsburgh ([@Piazza](#))** 



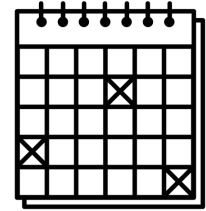
Module	Open time	Deadline
Phase 1 	2017-02-27 00:00	2017-03-26 23:59
Phase 2	2017-03-27 00:00	2017-04-16 16:59

Team Project Architecture



- Writeup and Queries will be released on Monday, **February 27th, 2017**
- We can have more discussions in subsequent recitations
- For now, ensure 3-person teams you decide have experience with web frameworks and database, storage principles and infra setup/hacking

Team Project Time Table



Phase (and query due)	Start	Deadline	Code and Report Due
Phase 1 • Q1, Q2, Q3	Monday 02/27/2017 00:00:00 ET	Sunday 03/26/2017 23:59:59 ET	Tuesday 03/28/2017 23:59:59 ET
Phase 2 • Q1, Q2, Q3, Q4	Monday 03/27/2017 00:00:00 ET	Sunday 04/16/2017 15:59:59 ET	
Phase 2 Live Test (Hbase/MySQL) • Q1, Q2, Q3, Q4	Sunday 04/16/2017 18:00:00 ET	Sunday 04/16/2017 23:59:59 ET	Tuesday 04/18/2017 23:59:59 ET
Phase 3 • Q1, Q2, Q3, Q4, Q5	Monday 04/17/2017 00:00:00 ET	Sunday 04/30/2017 15:59:59 ET	
Phase 3 Live Test • Q1, Q2, Q3, Q4, Q5	Sunday 04/30/2017 18:00:00 ET	Sunday 04/30/2017 23:59:59 ET	Tuesday 05/02/2017 23:59:59 ET

That's all for this morning!