## Course title: INFSCI 2750 Cloud Computing

**Description:** This course covers the fundamental concepts in Cloud Computing and provides hands-on experience in working with cloud systems. We will cover the introductory concepts of cloud software systems and provide an understanding of the design principles behind various existing cloud solutions. We will first cover the concepts of large scale parallel data processing in the cloud including the MapReduce programming model and Hadoop and its related ecosystem. We then focus on various existing virtualized commercial cloud models including the fundamental concepts of system virtualization, hypervisors and virtualized platforms. The next part of the course will focus on cloud storage systems including key-value stores and geographically distributed clouds. The last part of the course provides an introduction to security and privacy issues in cloud computing covering the issues of data and execution privacy and legal issues in modern commercial cloud services.

### Pre-requisites:

- TEL2000 OR Equivalent Background
- INFSCI 2500

# Tentative Schedule

| Date | Lecture topic | Objectives | Testing |
|---|---|---|---|
| Week 1 | Course Introduction<br><br>Introduction to Cloud Computing | *Describe/explain* the importance of elasticity, pay-per-use and the key distinctions between software-as-a-service, platform-as-a-service and infrastructure-as-a-service models.<br><br>Required reading: Reference [1] | |
| Week 2 | Introduction to Datacenters and Datacenter Systems<br><br>Google Filesystem (GFS) Case study | *Recognize/explain* the basic challenges in datacenter filessytems<br><br>*Identify* the key principles underlying | Quiz1 (based on required reading in week 1) |

| | | the design of large scale distributed filesystems<br><br>Required reading:[7] | |
|---|---|---|---|
| Week 3 | MapReduce Programming Model<br><br>MapReduce programming examples | *Recognize/explain* MapReduce programming model<br><br>*Explain and use* Piglatin and Hive commands for writing high-level language queries in Hadoop<br><br>Required reading: [9] | Mini Project 1 |
| Week 4 | Hadoop, Piglatin, Hive case studies<br><br>Overview of MapReduce optimization Techniques | *Recognize*, *explain* the inefficiencies in original Hadoop schedulers<br><br>*Understand, explain* various optimization techniques for in-memory MapReduce<br><br>Required Reading: [42] | Quiz 2 (based on required reading in weeks 2 and 3) |
| Week 5 | Introduction to virtualization and virtualized cloud platforms<br><br>Overview of Virtual Machine placement and live Virtual Machine Migration techniques | *Understand, explain* basic virtualization techniques in Cloud systems<br><br>*Recognize and analyze* various virtual machine placement and migration techniques<br><br>Required reading: [16] | |
| Week 6 | Cloud Storage systems and key | *Recognize/explain* the basic challenges in large scale storage | Quiz 3 |

| | | | |
|---|---|---|---|
| | value stores<br><br>Overview of Bigtable and Dynamo systems<br><br>Amazon Web Services | systems<br><br>*Identify* the key principles underlying the design of Bigtable and Dynamo systems<br><br>Required Reading:[8] | (based on required reading in weeks 4 and 5)<br><br>Mini Project 2 |
| Week7 | MapReduce in a Cloud: Challenges, Architectures and Techniques<br><br>Overview of Elastic MapReduce, Purlieus and Cura systems | *Explain and Recognize* basic challenges in dedicated MapReduce clouds<br><br>*Understand* the key design principles in existing MapReduce Cloud systems<br><br>Required Reading: [43] | |
| Week 8 | Mid term | | |
| Week 9 | Geographically distributed Clouds and overview of Spanner system | *Recognize, compare/contrast,* challenges of geographically distributed datacenter management with conventional datacenter resource management problem<br><br>*Recognize, identify* the design features of Spanner system | Quiz 4<br><br>(based on required reading in weeks 6 and 7) |

| | | Required reading: [20] | |
|---|---|---|---|
| Week 10 | Introduction to Blockchains and distributed consensus and ledger management . | *Recognize, identify* challenges of distributed consensus in blockchains<br><br>*Recognize, identify* the design features of Blockchain-based systems<br><br>Required reading: [46] | Mini Project 3 |
| Week 11 | Introduction to Security issues in Cloud Computing | *Understand, explain* security challenges in cloud computing<br><br>*Recognize*, *explain* the design principles of CryptDB system<br><br>Required reading: [45] | Quiz 5<br><br>(based on required reading in weeks 9 and 10) |
| Week 12 | Data Privacy in Cloud computing | *Recognize, explain* the basic privacy and legal issues in cloud systems<br><br>*Understand, explain* k- anonymity and differential privacy models and techniques | |

| | | Required reading: [19] | |
|---|---|---|---|
| Week 13 | Final Exam | | |

## Grading Policy

Reading Assignment-based Quiz (10%)

Three mini projects  (40%)

Midterm (20%)

Final Exam (20%)

Mini presentation (10%)

Class participation & Discussion (3% extra credit)

## References
1. "Above the Clouds: A Berkeley View of Cloud Computing", Michael Armbrust, et al. Technical Report, University of Berkerley, 2009, http://www.eecs.berkeley.edu/Pubs/TechRpts/2009/EECS-2009-28.pdf
2. "The Claremont Report on Database Research", 2008,
http://db.cs.berkeley.edu/claremont/claremontreport08.pdf
3. Hadoop, http://hadoop.apache.org/
4. Pig http://hadoop.apache.org/pig/
5. Hbase http://hadoop.apache.org/hbase/
6. Hive http://hadoop.apache.org/hive/
7. "The Google File System", Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung, OSDI, 2003, http://labs.google.com/papers/gfs-sosp2003.pdf
8. "Bigtable: A Distributed Storage System for Structured Data", Fay Chang, et al. OSDI 2006, http://labs.google.com/papers/bigtable-osdi06.pdf
9. "MapReduce: Simplified Data Processing on Large Clusters", Jeffrey Dean and Sanjay Ghemawat, OSDI 2004, http://labs.google.com/papers/mapreduce-osdi04.pdf
11. A comparison of approaches to large-scale data analysis. A. Pavlo et al. SIGMOD2009, http://database.cs.brown.edu/sigmod09/benchmarks-sigmod09.pdf
12. Amazon Web Services, http://aws.amazon.com/
13. Eucalyptus (http://www.eucalyptus.com/)
14. AppEngine http://code.google.com/appengine/
15. Azure http://www.microsoft.com/azure/
16. "Xen and the Art of Virtualization", Paul Barham, et al., SOSP 2003,
http://www.cl.cam.ac.uk/research/srg/netos/papers/2003-xensosp.pdf

17. "Benchmarking cloud serving systems with YCSB" Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, Russell Sears, ACM Symposium on Cloud Computing, 2010.

18. Cloud Security and Privacy: An Enterprise Perspective on Risks and Compliance (Theory in Practice) by Tim Mather, Subra

19. "Airavat: Security and Privacy for MapReduce ", Indrajit Roy, Srinath T.V. Setty, Ann Kilzer, Vitaly Shmatikov, Emmett Witchel. NSDI 2010

20. Corbett, James C; Dean, Jeffrey; Epstein, Michael; Fikes, Andrew; Frost, Christopher; Furman, JJ; Ghemawat, Sanjay; Gubarev, Andrey; Heiser, Christopher; Hochschild, Peter; Hsieh, Wilson; Kanthak, Sebastian; Kogan, Eugene; Li, Hongyi; Lloyd, Alexander; Melnik, Sergey; Mwaura, David; Nagle, David; Quinlan, Sean; Rao, Rajesh; Rolig, Lindsay; Saito, Yasushi; Szymaniak, Michal; Taylor, Christopher; Wang, Ruth; Woodford, Dale, "Spanner: Google's Globally-Distributed Database", *Proceedings of OSDI 2012* (Google), retrieved 18 September 2012.

21. Raluca Ada Popa, Catherine M. S. Redfield, Nickolai Zeldovich, and Hari Balakrishnan. CryptDB: Protecting Confidentiality with Encrypted Query Processing. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP)*, Cascais, Portugal, October 2011.

22. T. Do, T. Harter, Y. Liu, H. Gunawi, A. Arpaci-Dusseau, R. Arpaci-Dusseau, "HARDFS: Hardening HDFS with Selective and Lightweight Versioning", 11th USENIX Conference on File and Storage Technologies (FAST '13)

23. Paper: A Tale of Two Erasure Codes in HDFS, Mingyuan Xia, McGill University; Mohit Saxena, Mario Blaum, and David A. Pease, IBM Research Almaden, FAST 2015

24. Comet: An Active Distributed Key-Value Store, R. Geambasu et al, OSDI 2010

25. Paxos Made Transparent, Heming Cui, et al, SOSP 2015

26. Tango: distributed data structures over a shared log, M. Balakrishnan, et al, SOSP 2013

27. Succinct: Enabling Queries on Compressed Data, Rachit Agarwal, Anurag Khandelwal, and Ion Stoica, University of California, Berkeley, NSDI 2015

28. DryadLINQ: A System for General-Purpose Distributed Data-Parallel Computing Using a High-Level Language, Yuan Yu et al, OSDI 2008

29. Large-scale Incremental Processing Using Distributed Transactions and Notifications, D. Peng et al, OSDI 2010

30. Map-reduce-merge: simplified relational data processing on large clusters, H.-C. Yang et al, SIGMOD 2007

31. MapReduce Online, T. Condie et al, NSDI 2010

32. Building Consistent Transactions with Inconsistent Replication, Irene Zhang, Naveen Kr. Sharma, Adriana Szekeres, Arvind Krishnamurthy, Dan R. K. Ports, SOSP 2015

33. Yesquel: Scalable SQL storage for Web applications, Marcos K. Aguilera, et al, SOSP 20153

34. Salt: Combining ACID and BASE in a Distributed Database, Chao Xie, Chunzhi Su, Manos Kapritsos, Yang Wang, Navid Yaghmazadeh, Lorenzo Alvisi, and Prince Mahajan, OSDI 2014

35. Transaction chains: Achieving Serializability with Low Latency in Geo-Distributed Storage Systems, Y, Zhang et al, SOSP 2013

36. Extracting More Concurrency from Distributed Transactions, Shuai Mu, Yang Cui, Yang Zhang, Wyatt Lloyd, Jinyang Li, OSDI 2014

37. No compromises: distributed transactions with consistency, availability, and performance, Aleksandar Dragojević, Dushyanth Narayanan, Edmund B Nightingale, Matthew Renzelmann, Alex Shamis, Anirudh Badam, Miguel Castro, SOSP 2015

38. CalvinFS: Consistent WAN Replication and Scalable Metadata Management for Distributed File Systems, Alexander Thomson, Google; Daniel J. Abadi, Yale University, Usenix FAST 2015

39. Wormhole: Reliable Pub-Sub to Support Geo-replicated Internet Services, Yogeshwer Sharma, et al, Facebook, NSDI 2015

40. CubicRing: Enabling One-Hop Failure Detection and Recovery for Distributed In-Memory Storage Systems, Yiming Zhang, et al, NSDI 2015

41. BlinkDB: Queries with Bounded Errors and Bounded Response Times on Very Large Data, Sameer

Agarwal et al, Eurosys 2013[ ][ ][ ][ ]

42. Pig Latin: A Not-So-Foreign Language for Data Processing, SIGMOD 2008

43. Balaji Palanisamy, Aameek Singh, Ling Liu and Bhushan Jain, "Purlieus: Locality-aware Resource Allocation for MapReduce in a Cloud", Proc. of *IEEE/ACM Supercomputing (SC' 11)*

*44. COMET: Code Offload by Migrating Execution Transparently, OSDI 2012*

45. CryptDB: Protecting Confidentiality with Encrypted Query Processing, SOSP 2011

46. Changyu Dong et al. Betrayal, distrust, and rationality: Smart counter-collusion contracts for verifiable cloud computing. ACM CCS, 2017