**Datacenter Architecture**
**ECE 561**

**Meeting Time & Location**
W 1:25-3:55PM
Fitzpatrick Schiciano A

**Faculty**
Professor Benjamin Lee
Office Hours: 210 Hudson Hall, By Appointment
benjamin.c.lee@duke.edu

**Webpage.** http://people.duke.edu/~bcl15/class/class_ece561spr20.html

**Synopsis.** This course covers advanced topics in data centers with an emphasis on computer architecture and systems. This course surveys recent advances in processor, memory, network, and storage. And it surveys modern software systems in computing clouds. Discussion-oriented classes focus on in-depth analysis of readings. Students will learn to reason about datacenter performance and energy efficiency. Students will complete a collaborative research project.

After completing this course, students should be able to
- Understand datacenter hardware architectures and system software to manage them.
- Read architecture and systems papers critically.
- Write constructive paper reviews
- Identify open research problems in datacenter architecture
- Design and execute a research project to address an open research problem

Final project and paper required. Appropriate for graduate and advanced undergraduate students.

**Prerequisites.** ECE/CS 250 or equivalent required.

**Grading.** Participation/Discussion: 25%; Response Papers: 25%; Project/Paper: 50%

**Academic Policy.** University policy will be strictly enforced. Zero tolerance for cheating and/or plagiarism.

**Participation/Discussion.** This course uses a seminar, not a lecture, format. Each class covers particular topics from assigned papers. Students are expected to read the assigned papers and to prepare for course discussions. A student will be assigned to lead the discussion for each paper.

**Response Papers.** The students should prepare an insightful critique of the assigned papers due at the beginning of class. These response papers should take the form of a constructive paper review, including (1) summary, (2) strengths, (3) weaknesses, (4) directions for future work. These response papers should be no longer than one page per class. Papers will be evaluated for brevity and depth of insight. Submit response papers to Sakai.

**Project/Paper.** The course ends with a research project. Intermediate deliverables include a research statement, research plan, extended abstract, final paper, and oral presentation.

**Readings**

Jan 15      Syllabus and Administrivia
                  Benjamin Lee – Research Talk
                  Guided Discussion

Jan 22      **[Surveys]**
- Barroso, et al. "The datacenter as a computer: An introduction of the design of warehouse-scale machines," Synthesis Lecture, 2013
- Hazelwood et al. "Applied machine learning at Facebook: A datacenter infrastructure perspective," HPCA 2018

Jan 29      [**Software Systems I**]
- Ghemawat et al. "The Google file system," SOSP 2003.
- Dean and Ghemawat. "MapReduce: Simplified data processing on large clusters," OSDI 2004.
- Zaharia et al. "Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing." NSDI 2012.

Feb 5      **[Software Systems II]**
- Zaharia et al. "Discretized streams: Fault-tolerant streaming computation at scale." SOSP 2013.
- TensorFlow: A system for large-scale machine learning," OSDI 2016
- Li et al. "Scaling distributed machine learning with the parameter server," OSDI 2014.

Feb 12      **[Software Systems III]**
- Nakamoto. "Bitcoin: A peer-to-peer electronic cash system."
- Ethereum: A next-generation smart contract and decentralized application platform."
- Kanev et al. "Profiling a warehouse-scale computer," ISCA 2015.

Feb 19      **[Performance Analysis]**
- Dean. "The tail at scale" CACM 2013.
- Zheng and Lee. "Hound: Causal learning for datacenter-scale straggler diagnosis." SIGMETRICS 2018.
- Alipourfard et al. "CherryPick: Adaptively unearthing the best cloud configurations." NSDI 2017

Feb 26      **[Processors I]**
- Reddi et al., "Web search using mobile cores: Quantifying and mitigating the price of efficiency," ISCA 2010.
- Sze et al. "Efficient processing of deep neural networks: A tutorial and survey," Proceedings IEEE 2017.
- Jouppi et al. "In-datacenter performance analysis of a tensor processing unit," ISCA 2017.

Mar 4      **[Processors II]**
- Chung et al. "Single-chip heterogeneous computing: Does the future include custom logic, FPGAs, and GPGPUs?" MICRO 2010.

- Putnam et al. "A reconfigurable fabric for accelerating large-scale datacenter services," ISCA 2014.
- Kawaja et al. "Sharing, protection, and compatibility for reconfigurable fabric with AmorphOS," OSDI 2018.

Mar 18    **[Memory]**
- Malladi et al. "Towards energy-proportional datacenter memory with mobile DRAM," ISCA 2012.
- Ousterhout et al. "The case for RAMClouds," ACM SIGOPS Operating Systems Review, 2010.
- Nishtala et al. "Scaling Memcache at Facebook," NSDI 2013.

Mar 25    **[Network]**
- Al-Fares et al. "A scalable, commodity data center network architecture," SIGCOMM 2008.
- Benson et al. "Network traffic characteristics of data centers in the wild," SIGCOMM 2009.
- Heller et al. "ElasticTree: Saving energy in datacenter networks," NSDI 2010.

Apr 1    **[Resource Management]**
- Hindman et al. "Mesos: A platform for fine-grained resource sharing in the data center," NSDI 2011.
- Boutin et al. "Apollo: Scalable and coordinated scheduling for cloud-scale computing." OSDI 2014.
- Verma et al. "Large-scale cluster management at Google with Borg." EuroSys 2015.

April 8    **[Power Management]**
- Fan et al., "Power provisioning for a warehouse-sized computer," ISCA 2007.
- Meisner et al. "PowerNap: Eliminating server idle power." ASPLOS 2009.
- Qureshi et al. "Cutting the electrical bill for Internet-scale systems," SIGCOMM 2009.

April 15    **[Scheduling]**
- Ghodsi et al. "Dominant resource fairness: Fair allocation of multiple resource types," NSDI 2011.
- Zahedi et al. "Resource elasticity fairness with sharing incentives for multiprocessors," ASPLOS 2014.
- Fan et al. "The computational sprinting game," ASPLOS 2016.