

# Midterm Review

*DS 5110/CS 5501: Big Data Systems*  
*Spring 2024*

Yue Cheng



# Midterm exam

- Wednesday, February 28, 3:30 pm – 6:30 pm
  - Open book, open notes
- Covering four topics from Lec 2 to Lec 5
  - CPU job scheduling policies
  - Caching policies
  - MapReduce + HDFS
  - Spark

# Midterm exam

- The exam sheet will be available on **gradescope** at 3:30 pm (you will receive entry code after the class)
- You should work directly on the PDF document
  - Or, you may print it and write on printed papers, make sure you scan it to PDF with **visible resolution**
  - If you choose to scan using a smartphone camera, make sure it **covers everything clearly** – unrecognizable photos will not be graded
- Submission closes at 7pm
  - If you choose to scan, make sure your printer & scanner are handy

# CPU job scheduling

- FIFO
  - How it works?
  - FIFO's problems (why we need SJF)?
- SJF
  - How it works?
  - Any limitations (why we need STCF)?
- STCF (preemptive SJF)
  - How it works? How it solves SJF's limitations?
- RR (Round Robin)
  - How it works?

# CPU scheduling worksheet

# Caching policy

- LRU (least recently used)
- FIFO (first-in, first-out)

# MapReduce + HDFS

- How MapReduce works
- The performance characteristics of different phases of a MapReduce job (TeraSort)
- Fault tolerance
  - Replication for HDFS
  - Backup tasks for MapReduce

# Spark

- Motivation
- Transformations and actions
  - Narrow vs. wide transformation
- PageRank example
  - How iterative PR algorithm works
  - Optimizations on baseline PageRank
    - Co-partitioning for communication-efficient join
    - Apply `.persist(StorageLevel.DISK_ONLY)` for fault tolerance



# Question types

- Multi-choice questions (40%)
- True or false questions (25%)
- Problem solving (35%)

**Good Luck!**