



# Lecture 18: Job Scheduling

Abhinav Bhatele, Department of Computer Science



UNIVERSITY OF  
MARYLAND

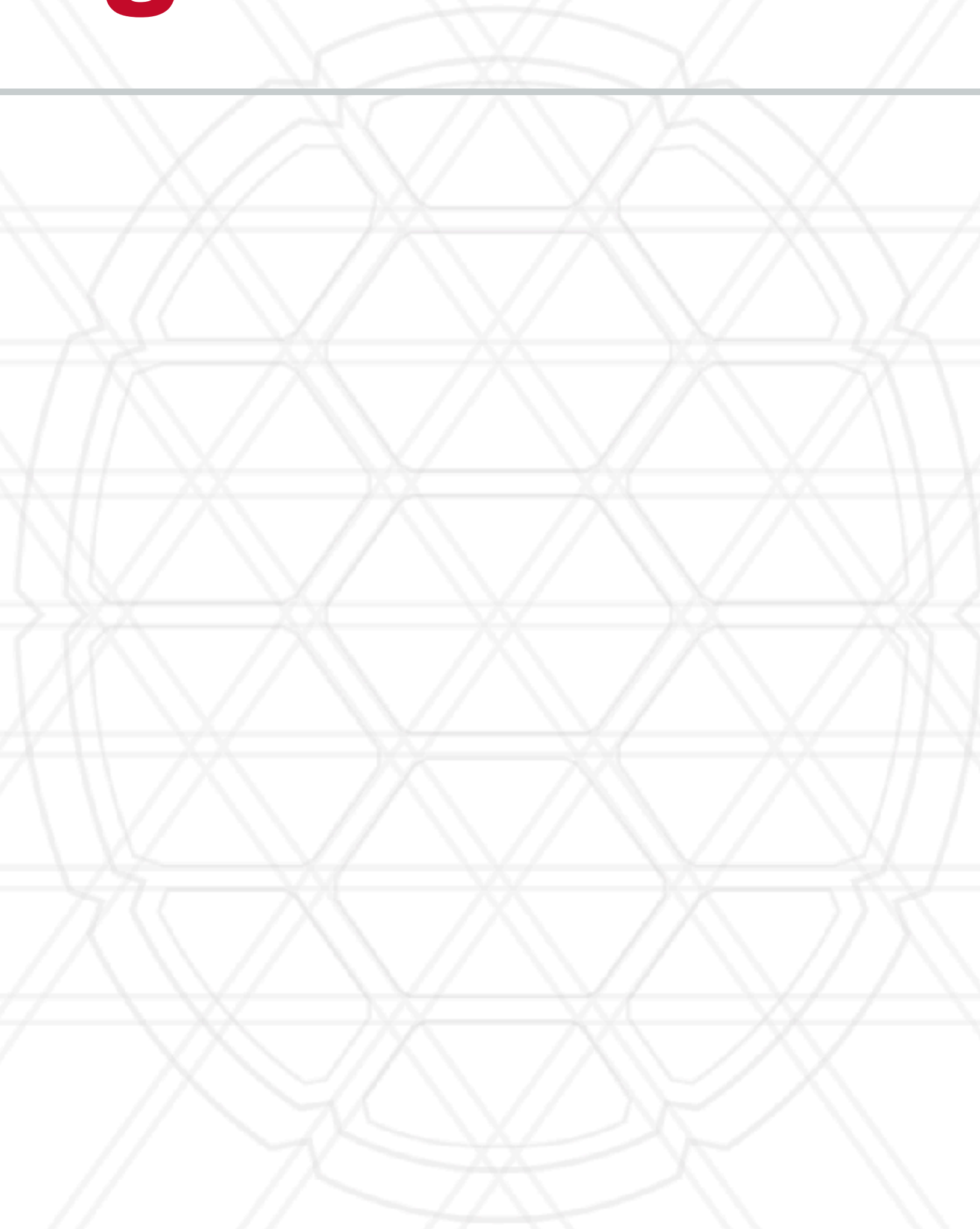
# Summary of last lecture

---

- OS daemons can lead to noise or jitter, which leads to performance variability
- Variability leads to practical issues and impacts software optimization cycle
- Can be mitigated by pinning processes and threads and leaving some cores free
- Can significantly impact performance of bulk synchronous programs
- Communication variability comes from other jobs sharing the same network

# Job scheduling

---


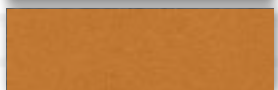

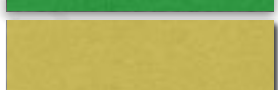

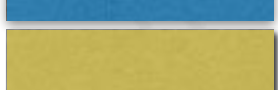


# Job scheduling

---

- HPC systems use job or batch scheduling
- Each user submits their parallel programs for execution to a “job” scheduler

## Job Queue

		#Nodes Requested	Time Requested
1		128	30 mins
2		64	24 hours
3		56	6 hours
4		192	12 hours
5		...	...
6		...	...

# Job scheduling

- HPC systems use job or batch scheduling
- Each user submits their parallel programs for execution to a “job” scheduler
- The scheduler decides:
  - what job to schedule next (based on an algorithm: FCFS, priority-based, ....)
  - what resources (compute nodes) to allocate to the ready job

## Job Queue

	#Nodes Requested	Time Requested
1	128	30 mins
2	64	24 hours
3	56	6 hours
4	192	12 hours
5	...	...
6	...	...

# Job scheduling

- HPC systems use job or batch scheduling
- Each user submits their parallel programs for execution to a “job” scheduler
- The scheduler decides:
  - what job to schedule next (based on an algorithm: FCFS, priority-based, ....)
  - what resources (compute nodes) to allocate to the ready job

- **Compute nodes: dedicated to each job**
- **Network, filesystem: shared by all jobs**

## Job Queue

	#Nodes Requested	Time Requested
1	128	30 mins
2	64	24 hours
3	56	6 hours
4	192	12 hours
5	...	...
6	...	...

# Job scheduling

- HPC systems use job or batch scheduling
- Each user submits their parallel programs for execution to a “job” scheduler
- The scheduler decides:
  - what job to schedule next (based on an algorithm: FCFS, priority-based, ....)
  - what resources (compute nodes) to allocate to the ready job

Concurrently running jobs can contend for shared resources: network, filesystem

- Compute nodes: dedicated to each job
- Network, filesystem: shared by all jobs

## Job Queue

	#Nodes Requested	Time Requested
1	128	30 mins
2	64	24 hours
3	56	6 hours
4	192	12 hours
5	...	...
6	...	...

# Two components of a scheduler

---

- Decide what job(s) to schedule next: scheduler
- Decide what nodes (and other resources) to allocate to them: resource manager



# Scheduling policies

---

- First come first serve (FCFS)
- Priority-based
  - Depending on project name and remaining allocation
- Backfilling
  - Use idle nodes that are being reserved for the next large jobs
  - Aggressive (EAZY) backfill: run jobs as long as they don't delay the first job (could lead to unbounded delays)
  - Conservative backfill: runs jobs as long as they don't delay **any** future job

# Resource management

---

- Most primitive: manage nodes
- Advanced management:
  - Node type aware (low vs. high memory, GPU nodes)
  - Network topology aware
  - Power aware

# Space sharing and time sharing

---

- Space sharing: Exclusive access to a resource until job completion
- Time sharing: Interleaved access to the same resource
  - Co-scheduling
  - Gang scheduling

# Quality of service metrics

---

- Job Wait Time: time between a job's submission and start

$$T_{\text{wait}} = T_{\text{start}} - T_{\text{submit}}$$

- Slowdown: incorporates running time of a job

$$\text{Slowdown} = \frac{T_{\text{wait}} + T_{\text{running}}}{T_{\text{running}}}$$

# Quality of service metrics

---

- System Utilization: fraction of nodes allocated to running jobs at a given time

$$utilization_t = \frac{N_t}{N}$$

- Schedule Makespan: time between the first job's submission and last job's completion for a job trace (workload)

# Questions?



UNIVERSITY OF  
MARYLAND

Abhinav Bhatele

5218 Brendan Iribe Center (IRB) / College Park, MD 20742

phone: 301.405.4507 / e-mail: [bhatele@cs.umd.edu](mailto:bhatele@cs.umd.edu)