

Scheduling Criteria

- **Throughput** – # of threads that complete per unit time
 $\# \text{ jobs/time}$ (Higher is better)
 - **Turnaround time** – time for each thread to complete
 $T_{\text{finish}} - T_{\text{start}}$ (Lower is better)
 - **Response time** – time from request to first response ()
i.e. time between waiting to ready transition and ready to running transition
 $T_{\text{response}} - T_{\text{request}}$ (Lower is better)
- ➔ Above criteria are affected by secondary criteria
- CPU utilization – %CPU fraction of time CPU doing productive work
 - Waiting time – $\text{Avg}(T_{\text{wait}})$ time each thread waits in the ready queue

How to balance criteria?

- **Batch systems** (supercomputers)
strive for job throughput and turnaround time
- **Interactive systems** (personal computers)
strive to minimize response time for interactive jobs

However, in practice, users prefer predictable response time over faster but highly variable response time

Often optimized for an average response time