

# Discrete Event Simulator

---

Anshul Gupta (16305R001)

# Features

- Multi-core server machine
- Multi-threaded server
- Round-robin scheduling
- Closed loop system

# Input

- n\_users
- n\_CPUs
- easeInTime
- maxIters
- bufferCapacity
- threadpoolSize
- quantum
- ctxSwitchTime
- serviceTimeMean
- requestTimeoutMin
- requestTimeoutMax
- thinkTimeMean
- thinkTimeStdv
- retryThinkTimeMean
- retryThinkTimeStdv

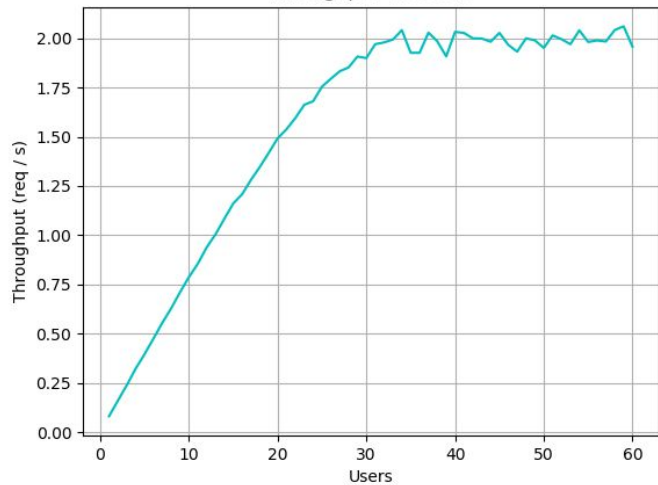
# Output

- arrivalRate
- throughput
- goodput
- badput
- responseTime
- utilization
- contextSwitchBusyTimeFraction
- requestsInSystem
- droppedFraction
- dropRate
- timedOutFraction

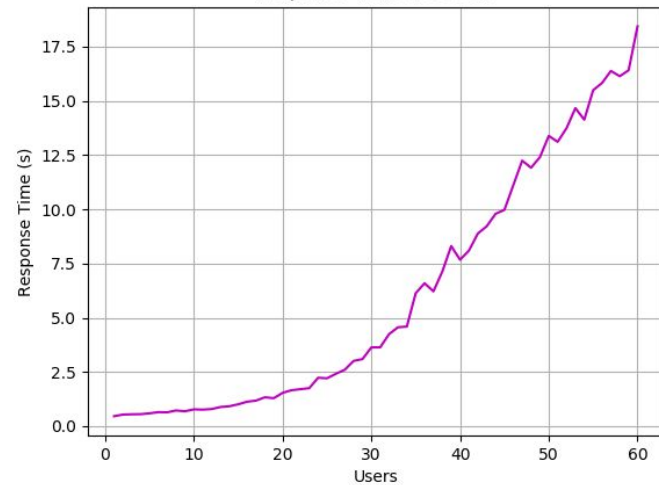
# Correctness of the code

- Results for M/M/1 queuing system are well known
- Increasing the quantum size of round-robin will simulate FIFO scheduling
- Setting buffer size to a large number will act as infinite buffer
- $\tau = 0.5$  sec,  $\mu = 2.0$  requests / sec
- Quantum size = 10000 sec
- Buffer Size = 8000000
- $n\_CPUs = 1$
- Think Time = 12 sec
- $M^* = 25$

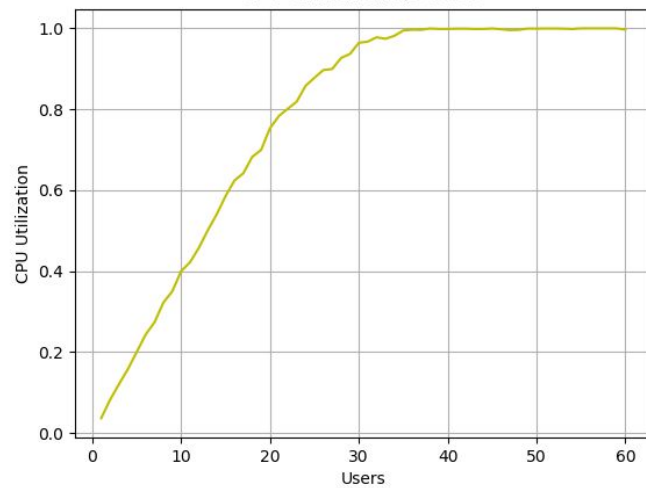
Throughput v/s Users



Response Time v/s Users



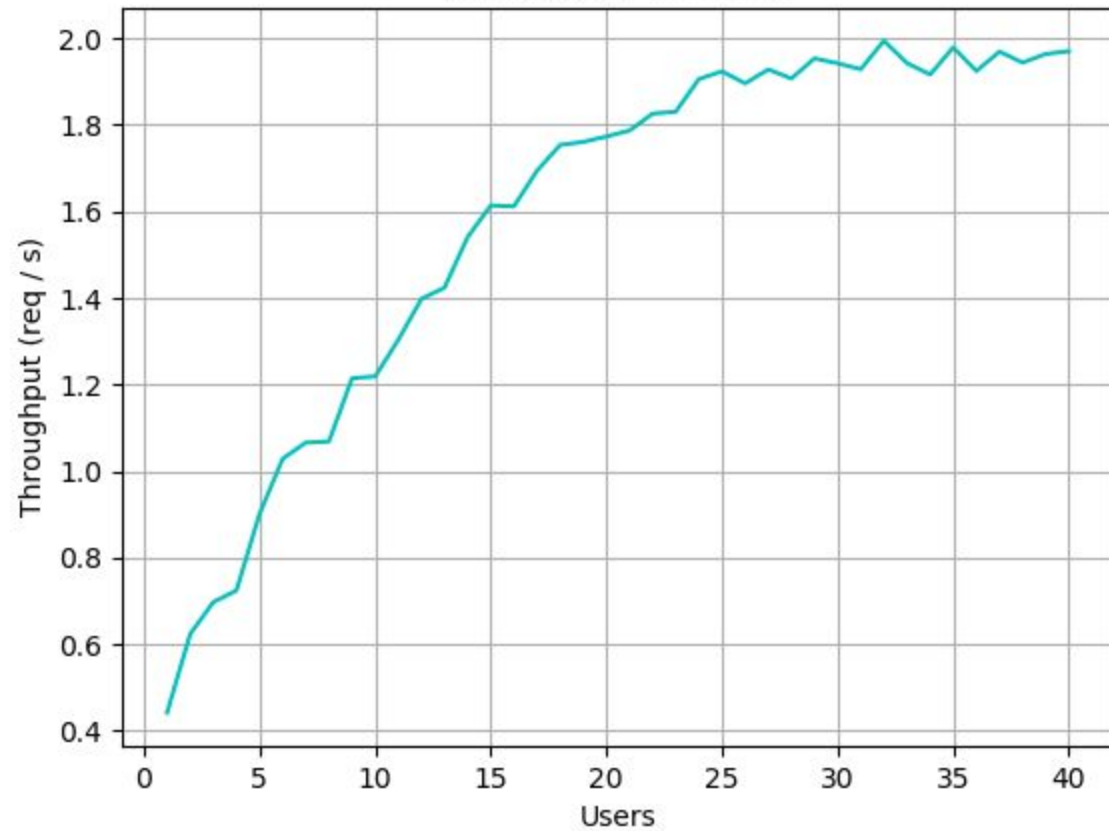
CPU Utilization v/s Users



# Simulation Parameters - 1 (Number of users)

- `n_users = [1, 40]`
- `n_CPUs = 4`
- `easeInTime = 20.0`
- `maxIters = 10000`
- `bufferCapacity = 800`
- `threadpoolSize = 200`
- `quantum = 0.5`
- `ctxSwitchTime = 0.01`
- `serviceTimeMean = 2`
- `requestTimeoutMin = 5.0`
- `requestTimeoutMax = 15.0`
- `thinkTimeMean = 12.0`
- `thinkTimeStdv = 4.0`
- `retryThinkTimeMean = 12.0`
- `retryThinkTimeStdv = 4.0`

Throughput v/s Users

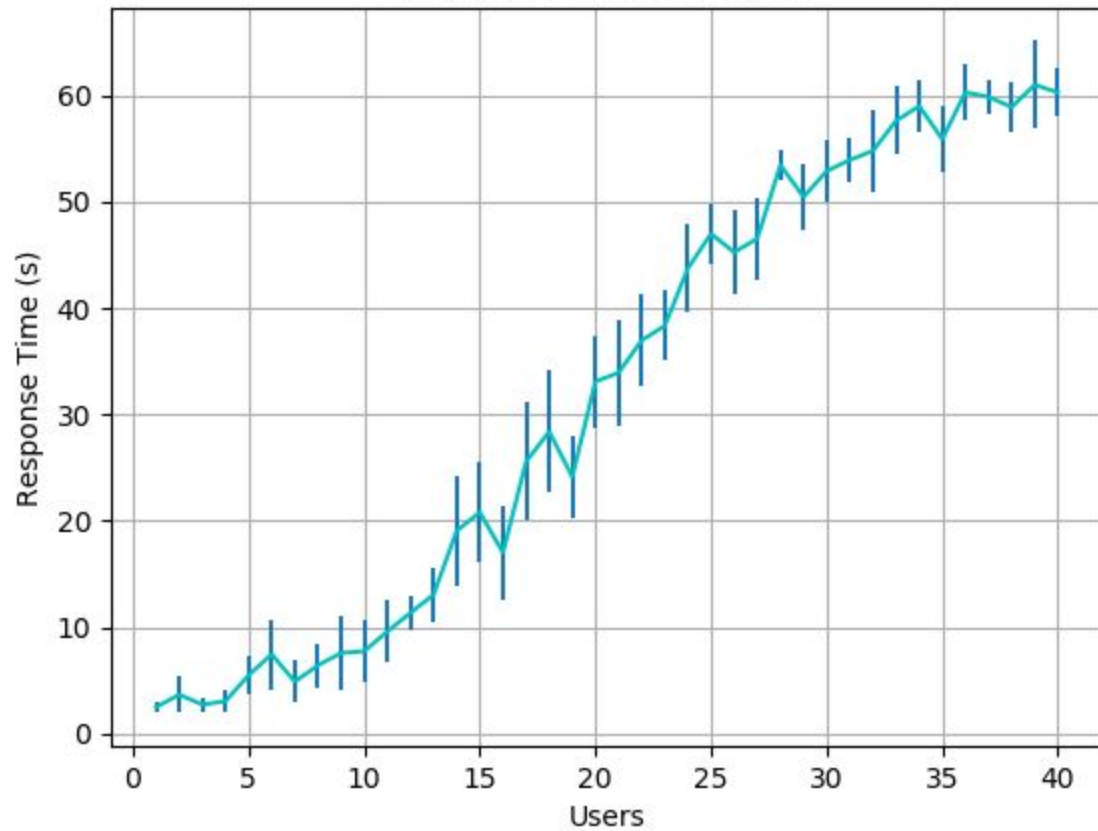




# Observations

- Until saturation, throughput increases almost linearly
- At high load, it saturates at 2 requests / sec

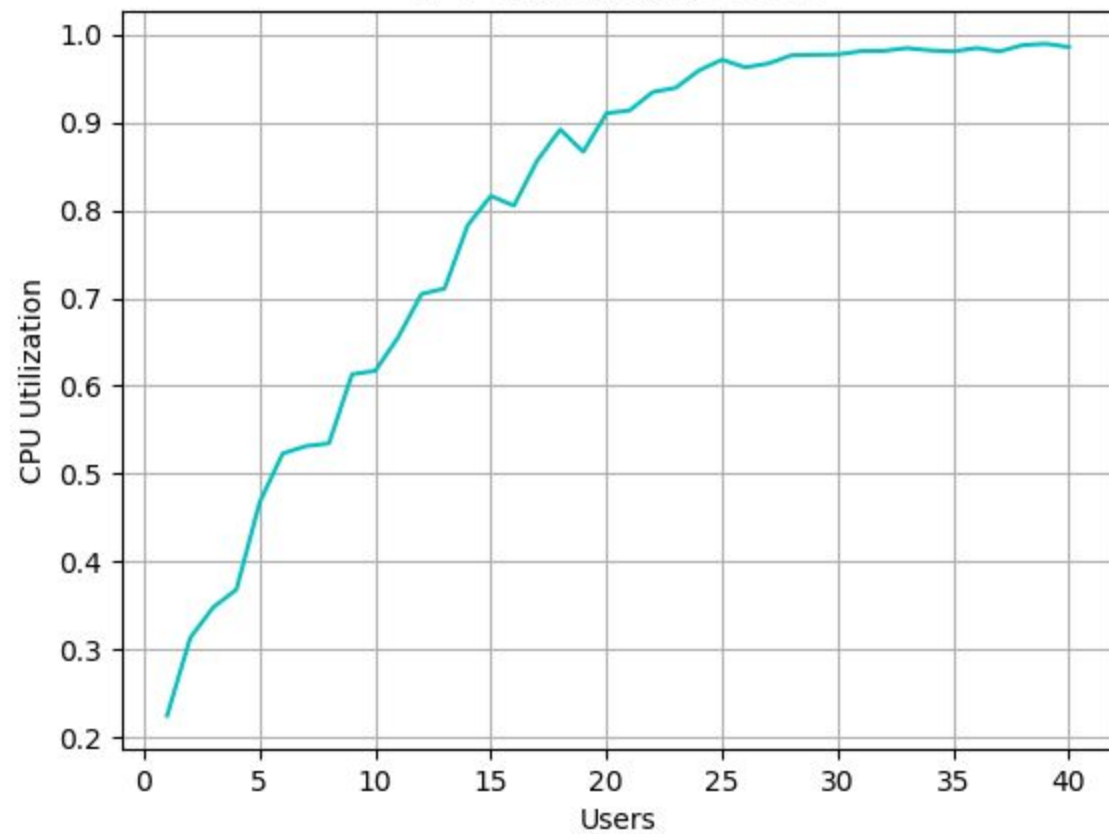
Response Time v/s Users



# Observations

- At low load, the response time is same as the service time (2 sec)
- After saturation point, response time increases sharply

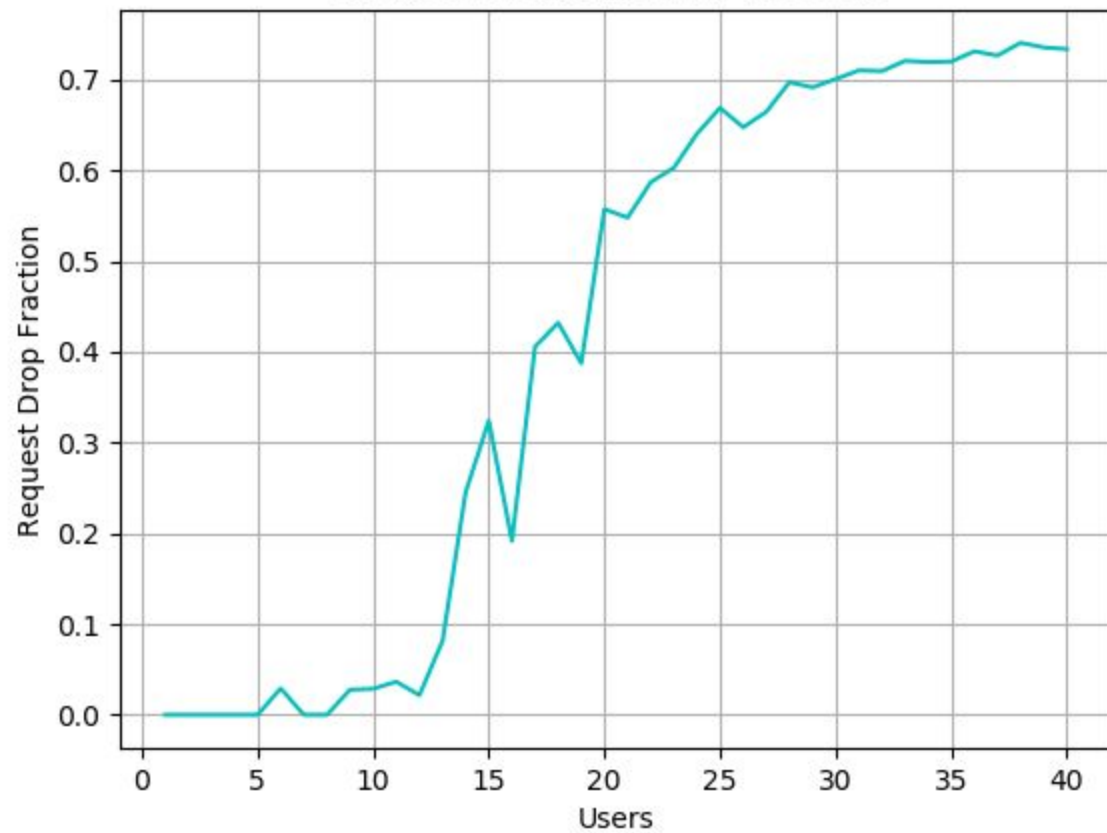
CPU Utilization v/s Users



# Observations

- As expected, the utilization is ~80% at  $M^* = 25$
- It saturates at around 30 - 32 users

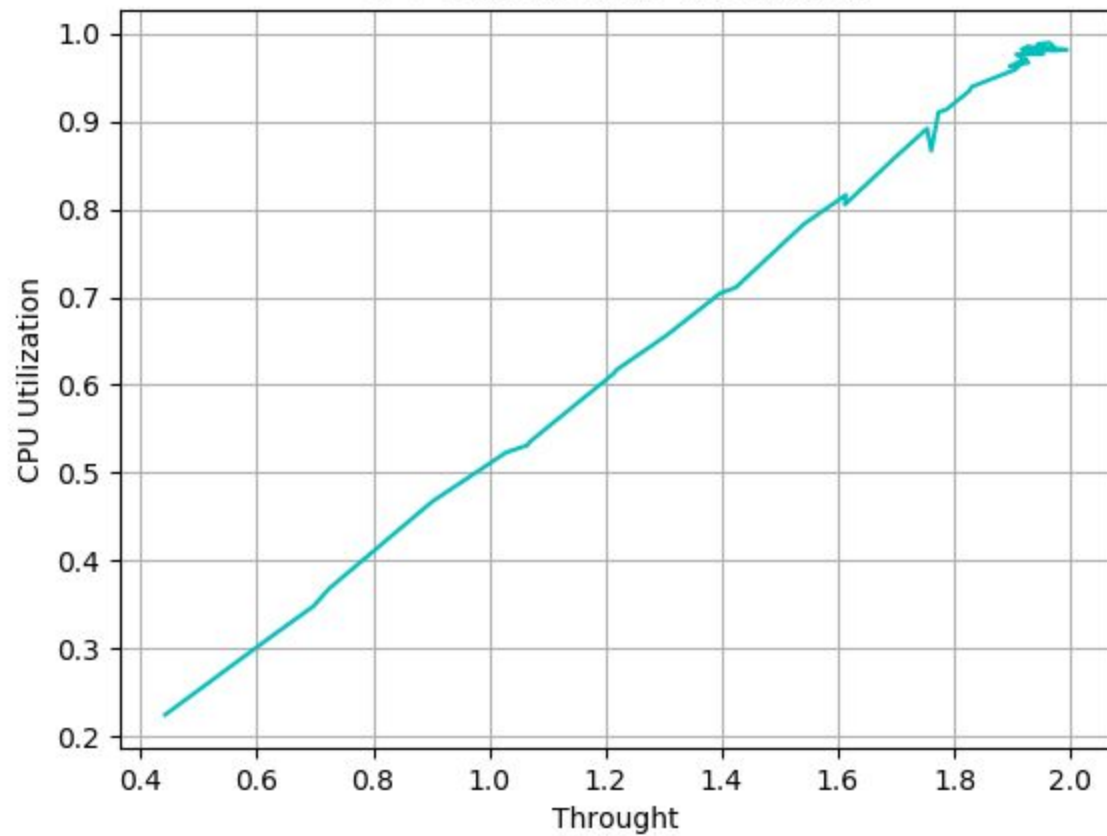
Request Drop Fraction v/s Users



# Observations

- As the number of users increases, the limited buffer cannot accommodate the requests and starts dropping them

CPU Utilization v/s Throughput





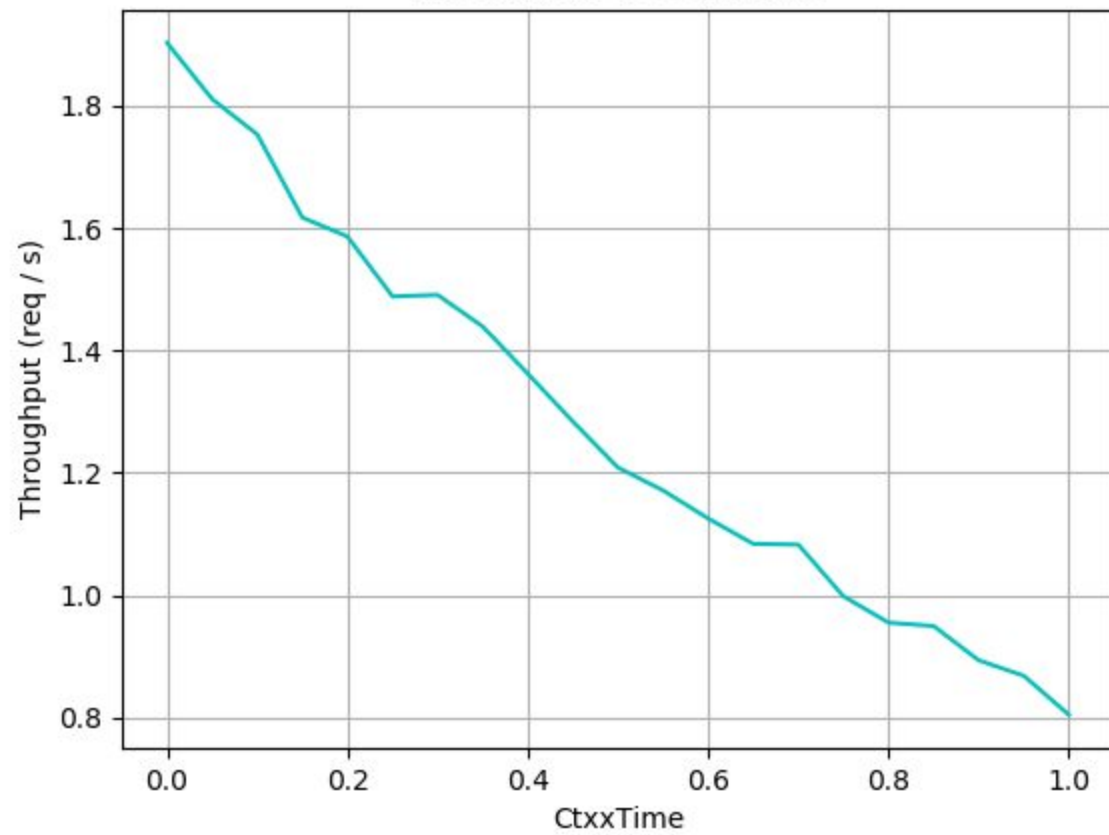
# Observations

- The linear relationship between throughput and utilization is maintained in the graph.
- Utilization is 1 when throughput saturates at 2 requests / sec

## Simulation Parameters - 2 (Context Switching Time)

- `n_users = 25`
- `n_CPUs = 4`
- `easeInTime = 20.0`
- `maxIters = 10000`
- `bufferCapacity = 800`
- `threadpoolSize = 200`
- `quantum = 0.5`
- `ctxSwitchTime = [0.0, 1.0, 0.05]`
- `serviceTimeMean = 2`
- `requestTimeoutMin = 5.0`
- `requestTimeoutMax = 15.0`
- `thinkTimeMean = 12.0`
- `thinkTimeStdv = 4.0`
- `retryThinkTimeMean = 12.0`
- `retryThinkTimeStdv = 4.0`

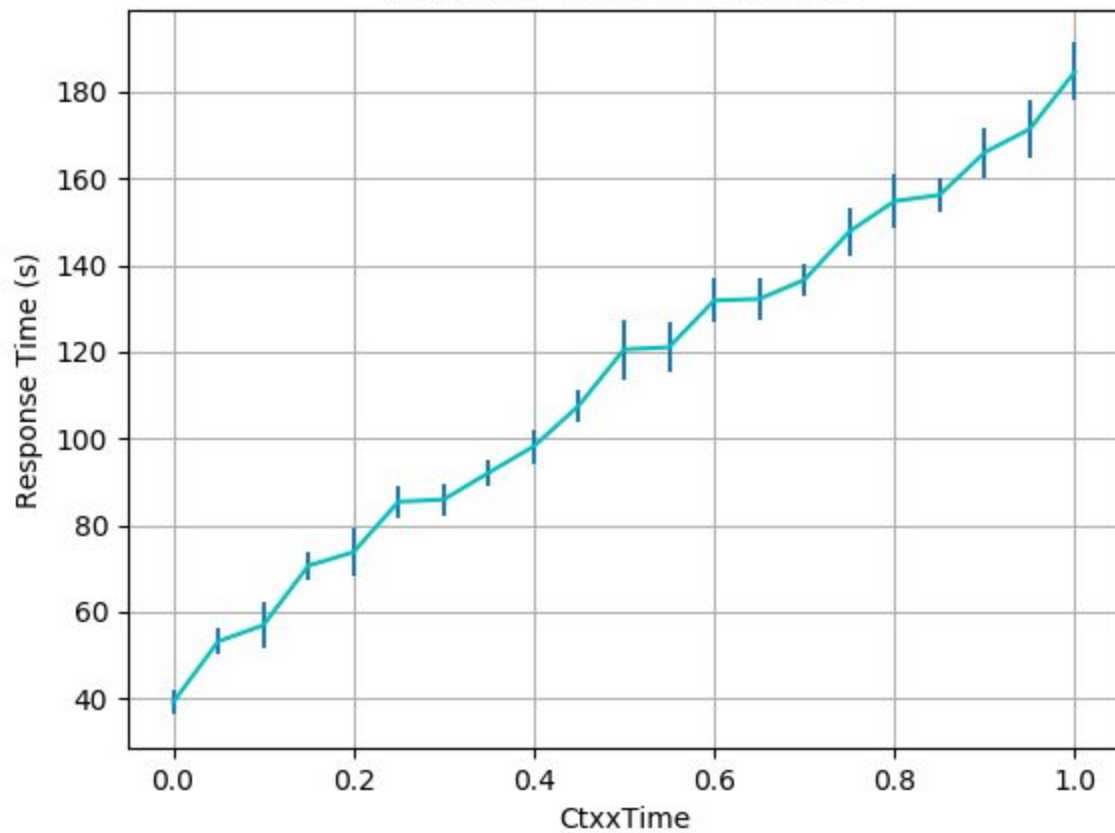
Throughput v/s CtxxTime



# Observations

- As the context switching time increases, it is obvious that the throughput will decrease
- Higher context switching time will cause system to reach saturation faster

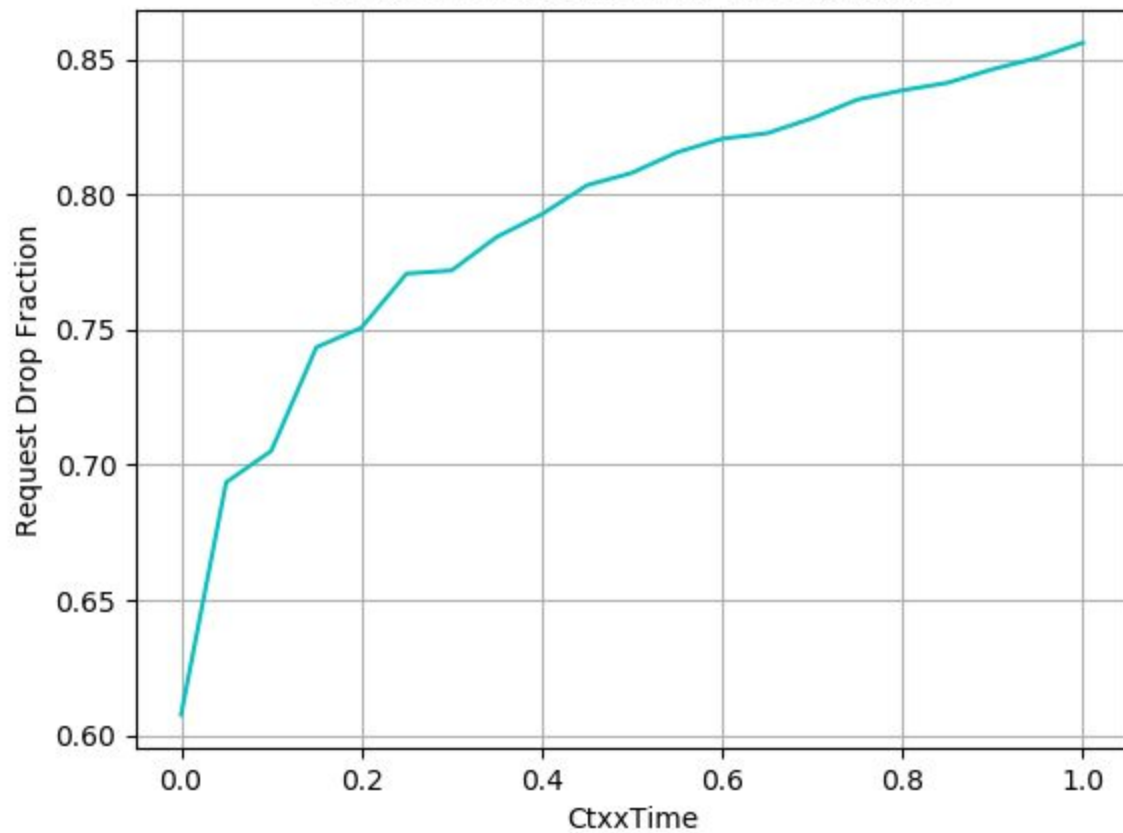
Response Time v/s CtxxTime



# Observations

- Similar reasoning as that of throughput
- Higher context switching time means more delays between processing the requests

Request Drop Fraction v/s CtxxTime



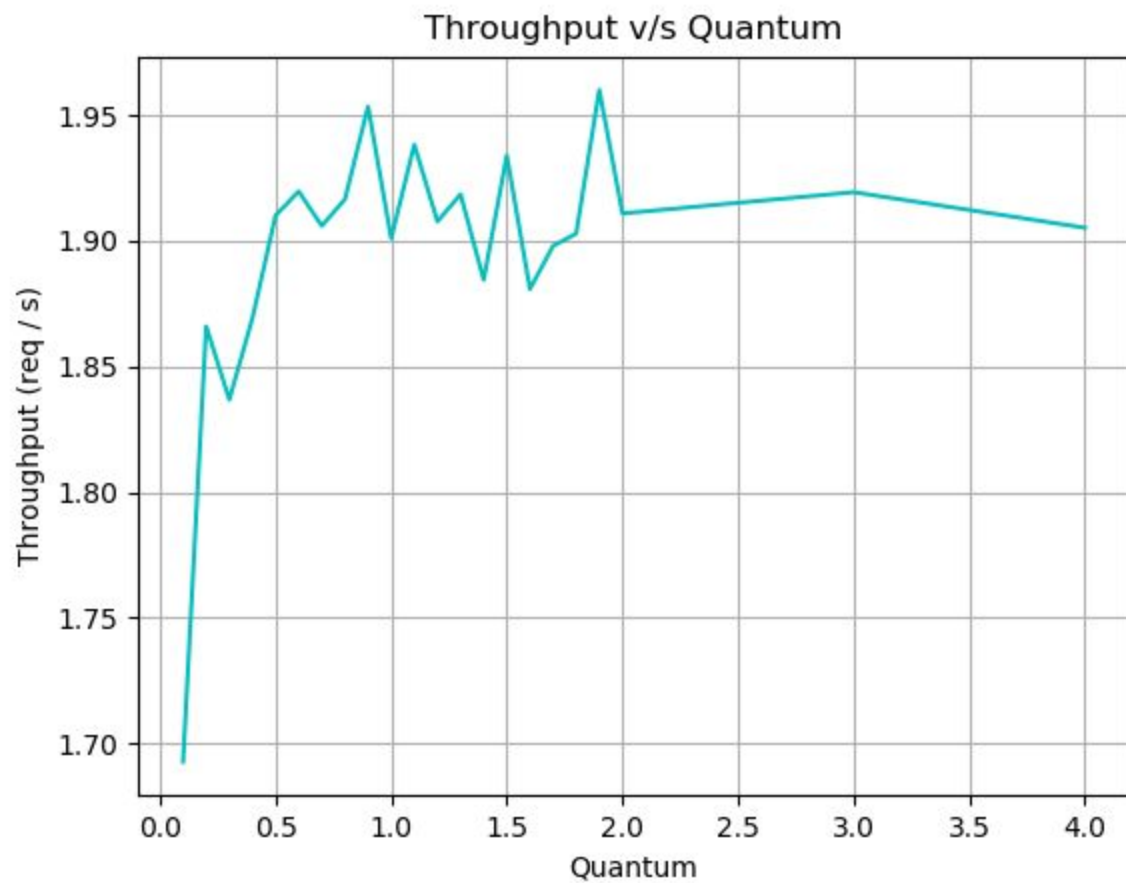
# Observations

- Since the requests take more time to process, the requests in buffer has to stay longer
- The buffer empties at a lower rate so the requests drop rate increases



# Simulation Parameters - 3 (Quantum)

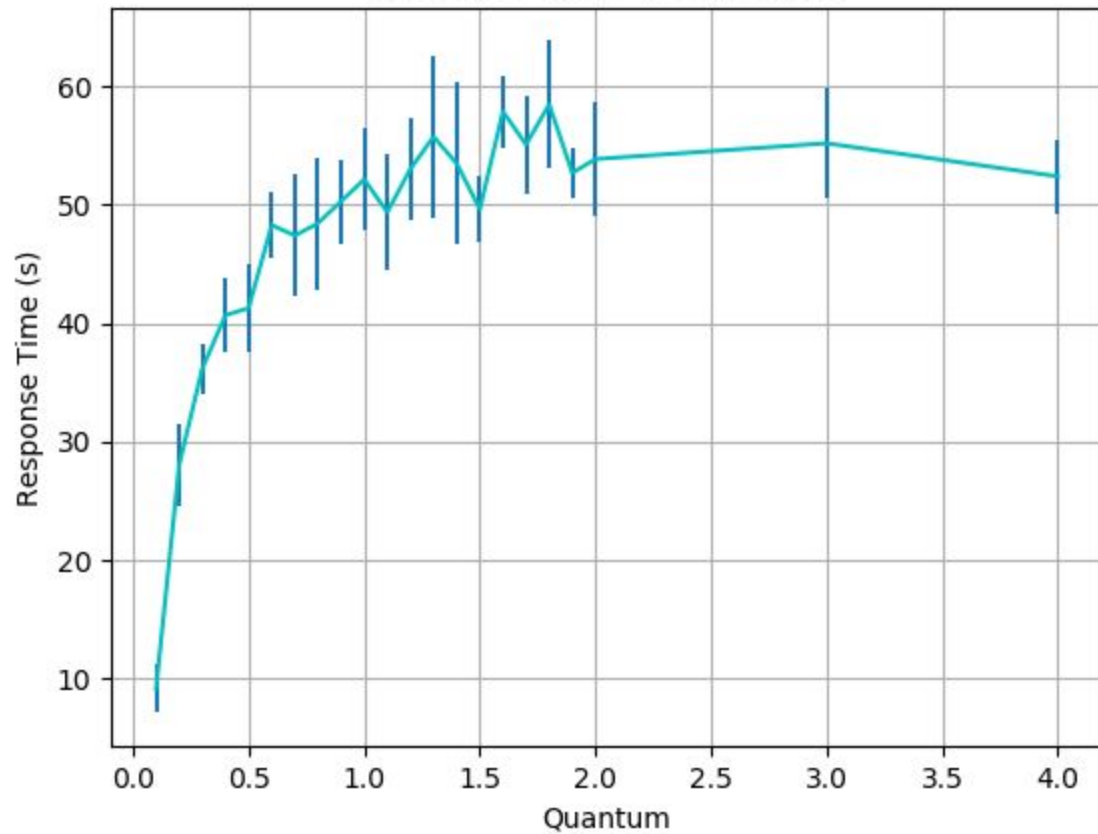
- `n_users` = 25
- `n_CPUs` = 4
- `easeInTime` = 20.0
- `maxIters` = 10000
- `bufferCapacity` = 800
- `threadpoolSize` = 200
- `quantum` = [0.1, 2.0, 0.1] + (3.0, 4.0)
- `ctxSwitchTime` = 0.01
- `serviceTimeMean` = 2
- `requestTimeoutMin` = 5.0
- `requestTimeoutMax` = 15.0
- `thinkTimeMean` = 12.0
- `thinkTimeStdv` = 4.0
- `retryThinkTimeMean` = 12.0
- `retryThinkTimeStdv` = 4.0



# Observations

- When quantum size is very low, there are more context switches therefore less throughput
- When quantum size is very high, there are very less context switches. So the throughput saturates as quantum is much higher than the service time of the request

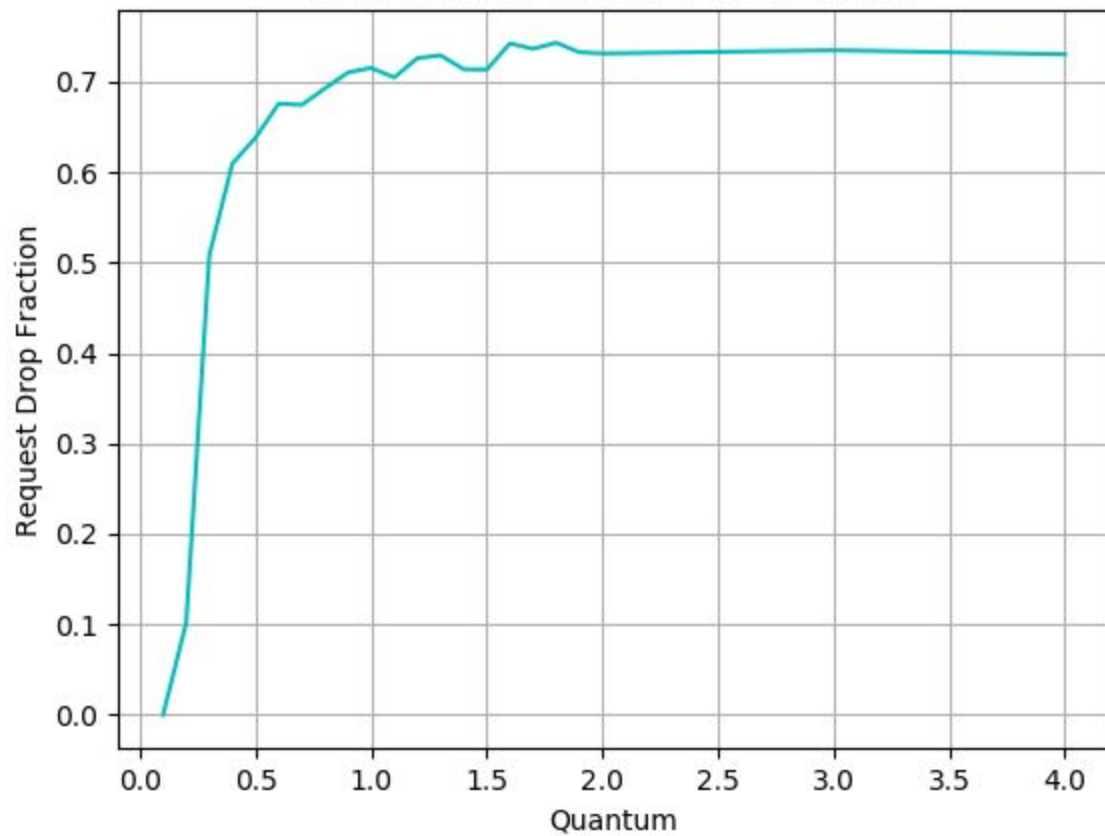
Response Time v/s Quantum



# Observations (Bug)

- Unable to explain the graph
- As the context switches are very high (low quantum), response time should be high and it should decrease as context switches decrease
- The graph is contradicting this

Request Drop Fraction v/s Quantum



# Observations (Bug)

- Since the throughput is low at low quantum, the requests stay in buffer for a longer time so the new requests should be dropped as the buffer is full
- Similarly at high quantum, the requests get processed very fast and the buffer clears at a higher rate. So requests drop rate should be low
- The graph is contradicting this

End

---