

CS681 Assignment 1 Part 1

Anand Dhoot (130070009)

Anchit Gupta (13D100032)

Problem Statement

Understand the performance of a queuing system, taking measurements and seeing how changes in parameters affect the metrics.

Specifically, we consider a web-server, with open-load and set the other control parameters such that we are able to observe the server CPU as the resource.

The open-load is generated using Httpperf and its output used to record all client side metrics

Home-brew script for measuring server CPU%

Experimental Setup

Server - Dell Inspiron - i5(one active core), 6GB RAM

Client - Lenovo Y500 - i7, 8GB RAM

Network - LAN (100 Mbps cables)

Httpperf - Hog mode, Timeout $+\infty$

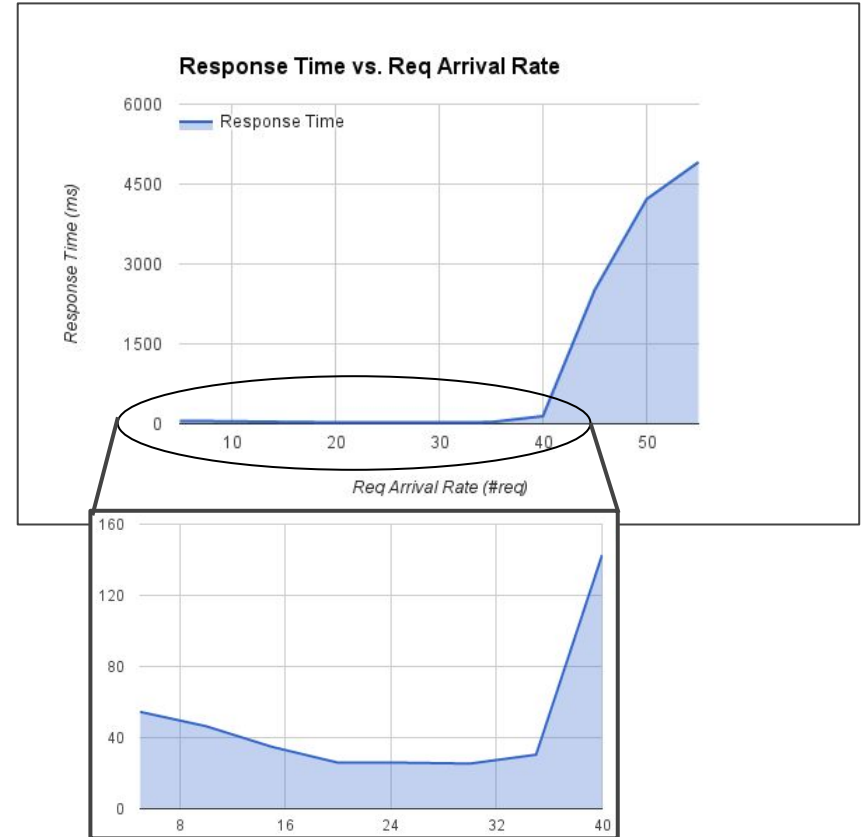
Theoretical Predictions

- Single-server queuing system (D/D/1)
- Predictions of metrics
 - Response Time - Constant then sharp linear increase
 - Server CPU Utilization - Linear increase and eventual saturation
 - Throughput - Linear increase and eventual saturation
- Utilization Law Prediction - Constant service time across the runs

Experimental Results

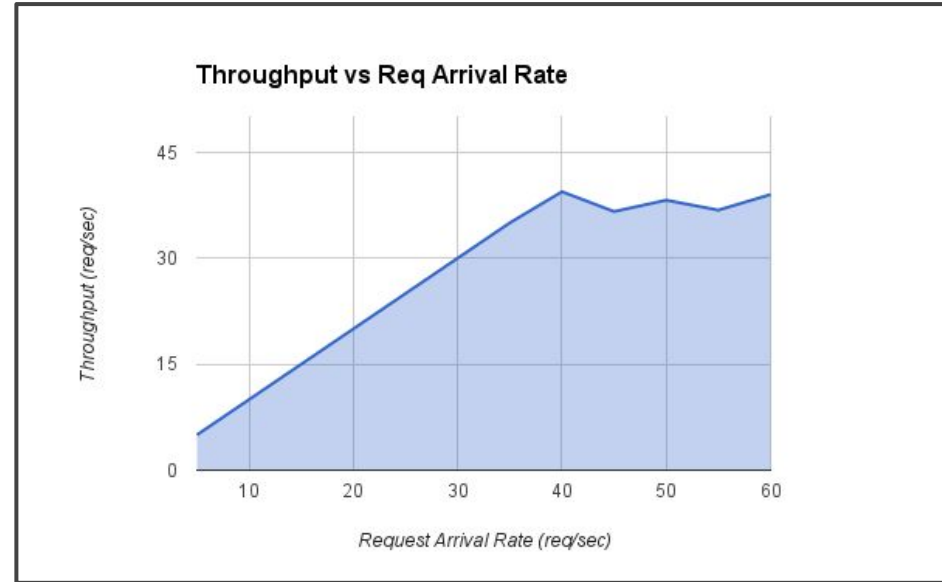
Response Time vs Request Arrival Rate

- Largely, as expected.
- Small dip in beginning - Cannot explain
- Constant for some time as queuing delay is negligible
- Sharp increase at ~40 req/sec. Subsequent increase linear as Queuing Delay becomes significant
- Increases unboundedly (no apparent saturation) as Timeout arbitrarily large



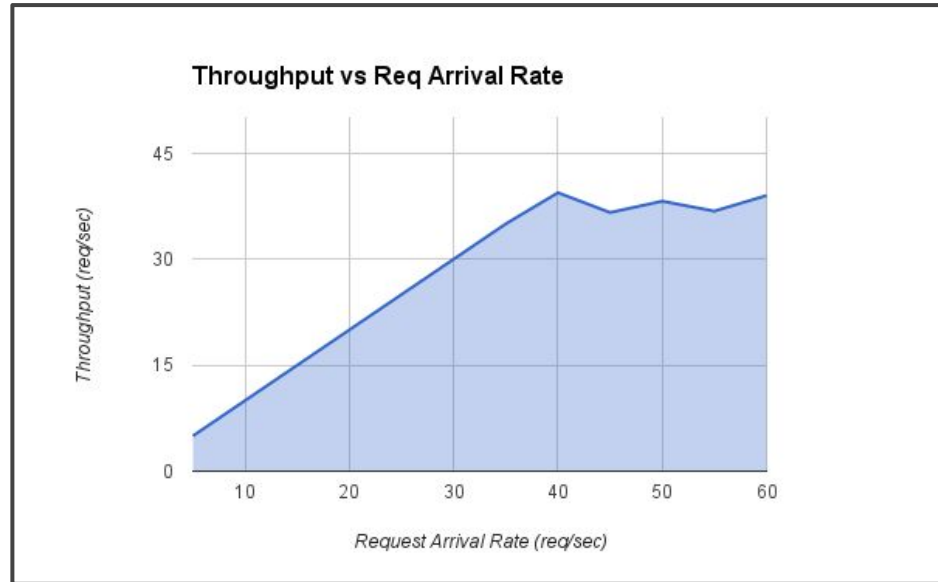
Throughput vs Request Arrival Rate

- Increases linearly initially - server able to process all requests
- Reached saturation when rate exceeds 40 req/sec
- Throughput fluctuates around 39 req/s for higher values of Request Arrival Rate



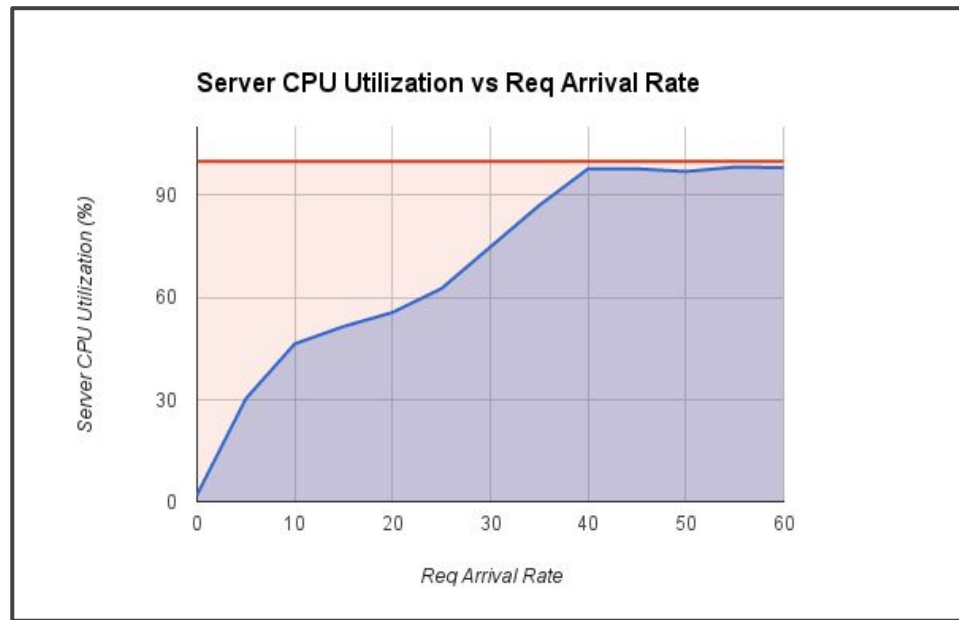
Throughput vs Request Arrival Rate

- This fluctuation is in line with the CPU% (they both seem to dip together)
- When the OS eats up a few extra CPU% this results in apache getting a bit less in that experiment and hence a little less throughput



(Server) CPU Utilization vs Request Arrival Rate

- Increasing initially. This is because server can use more of idle CPU to process more requests
- There is a small fluctuation near the saturation values (~1-2%)
- This is because we only a fraction of CPU is used by apache threads. In those runs the OS might be doing its own work resulting in a 1-2% fluctuation



Unsuccessful Requests vs Request Arrival Rate

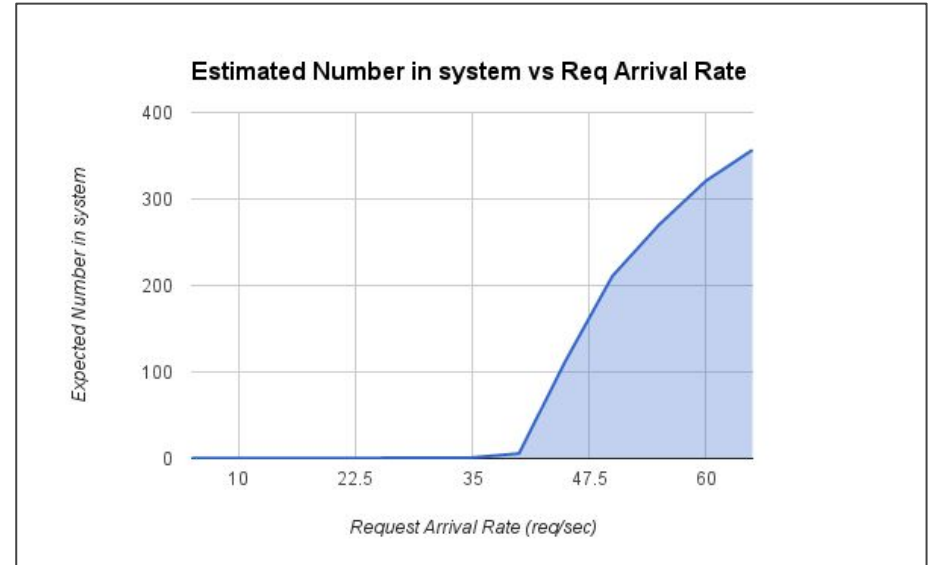
- No errors due to Connection issues (refused/rejected), Client/Socket timeout, FD Unavailable, etc.
- We ensured that our system had enough FD available, and the timeout was set arbitrarily large
- Zero throughput (as expected)

Estimating Service Time - Utilization Law

- We used the Utilisation Law to estimate the service time.
- I.e $\text{Service Time(ms)} = (\% \text{CPU} / \text{Throughput}) * 10$
- As expected this was roughly equal across the various runs
- We averaged it over the different experiments and found it out to be = 26.91 ms
- This corresponds to the server capacity of $\sim 1000 / 26.91 = 37.16$ req/sec, which is close to the experimentally measured estimate

Little's Law

- We estimated the no. of requests in the server using Little's Law.
- The effective arrival rate is same as the arrival rate (as number of threads is not a bottleneck)
- Graph similar in shape to response time. This is because major part of response time is queuing delay which is proportional to queue length



Conclusion

Server saturates at 39 req/sec

Average service time at 27 ms

Results seem reasonably consistent
with theoretical expectation