

Reducing Data Center Latency

Zainab Ghadiyali
zainab@cs.wisc.edu

Michael Griepentrog
mgriepentrog@cs.wisc.edu

Aditya Akella
akella@cs.wisc.edu

Abstract

Latency is usually compromised in order to improve bandwidth. However, with the advent of latency sensitive applications such as high frequency trading, the focus on reducing latency in network has increased. Furthermore, computation is increasingly also shifting to datacenters. Coupled together, these two factors necessitate further study on latency in datacenters. In this body of work, we study various factors that may influence latency, such as type of flows, offloading parameters, hardware such as network interface cards, optimized algorithms such as stochastic fair queueing and type of routers (presence and absence of switches). By utilizing a server-client architecture and OpenFlow switches, we show that offloading, stochastic fair queueing and utilization switches over direct flows help reduce latency in networks. Further, we quantify the difference and propose further experiments that would help identify other latency influencing factors.

1 Introduction

In the last few decades, the key focus of the network community has been on improving the overall goodput of networks. The initial focus on circuit switching moved to packet switching due to bandwidth and hardware inefficiencies in working with network resources for bursty communication traffic. The Transmission Control Protocol(TCP) was invented to address the issue of bandwidth/congestion collapse in the network and ensure bandwidth fairness. [1, 5, 7]

However network latency has been given less importance. The last 30 decades have brought a meager 30x reduction in latency, versus the 3000x times improvement in bandwidth. Several studies have been conducted in the 90s to tackle the latency issue, but they were not adopted. [3, 4, 8, 6] Additionally, most applications were throughput oriented (e.g. email) and hence not sensitive of delivery time. For example, switches are especially designed

with large packet buffers in order to avoid congestion since TCP protocol does not deal well with packet drops from congestion. Network Interface Cards (NICs) are optimized to delay interrupts for upto 30 μ s. Latency for latency sensitive applications is seen as an acceptable trade-off in order to maintain high bandwidth utilization. However, there is now an increasing interest in reducing latency in datacenters. A substantial amount of computing is now shifting to data centers and reducing latency is now easier given the confines of a building rather than tackling the issue for the internet at large. Furthermore, ultra low latency sensitive applications such as High Frequency Trading, HPC and Google Instant Search service could especially benefit from this. Several high-frequency traders implement trading applications with a goal to reach a trade decision in under 100 microseconds. NIC, end-host stacks and switches are all points where a packet may experience latency while traversing from client to server. In order to address delays arising from queueing delay at switches, DCTCP [1] uses ECN marking to slow down flows before the queue is saturated. HULL [2] further suggests trading off a little bandwidth to provide smaller amount of latency. In this paper, we look into understanding how switch, various offloading parameters and hardware influences latency.

OUR OBSERVATION IS THAT

THROUGH OUR EVALUATION WE FIND THAT

2 Methodology

2.1 Design

A strong evaluation and understanding of reducing latency will be incomplete without simultaneously studying the trade-offs between latency and other important factors such as bandwidth and CPU utilization. Thus, we study latency along with its influence on CPU utilization and bandwidth. Furthermore, it would be important to also understand the role of latency influenc-

ing factors in UDP and TCP. UDP should be faster than TCP since it allows a continuous packet stream versus TCP that sends acknowledgements (ACKS) for a set of packets calculated using the TCP window size and round trip time (RTT) Thus, metrics were obtained for TCP and UDP flows. Offloading moves the IP and TCP processing to the NICs. Keeping type of NIC constant, we can study the influence of offloading on latency. TCP segmentation offload (TSO)/ Generic segmentation offload (GSO) are considered useful in increase out-bound throughput of high-bandwidth network connections since it reduces host CPU cycles for protocol header processing and checksumming. Generic receiving offload (GRO) attempts to replicate the TSO modus operandi on the receiver-side. Effects of offloading is studied by:

1. TSO/GSO and GRO : On
2. TSO/GSO and GRO : Off
3. TSO/GSO : On and GRO : Off
4. TSO/GSO : Off and GRO : On

2.2 Experimental Setup

Two commodity computers connected via an OpenFlow switch [9]. One of the computers served as a client and the other as server. Eight experiments were ran in this setting. Another eight experiments ran on the two machines in absence of OpenFlow switch. For each experiment, throughput is recorded using iperf, cpu utilization through pidstat, interpacket delay through tcpdump and latency through ping. This process was automated by running bash scripts on server and client end. Offloading parameters were varied using ethtool. Data was analysed using wireshark and a python script that parsed results from experiments and provided basic summary in terms of min, max, stdev and average latency, average throughput in terms of Mbits/sec and average percent CPU utilization.

3 Evaluation

3.1 Bandwidth

TSO works by moving work of breaking down large chunks of data over to the NIC. In figavgBandwidth-ClientNoswitch, the x-axis represents offloading settings while y-axis represents bandwidth. Bandwidth is almost halved when TSO/GSO if switched off. This suggests that throughput may be increased by switching off offloading. Keeping GRO on in the absence of TSO/GSO did not show an appreciable change, thus suggested that GRO plays a little to no role in increasing bandwidth.

Interestingly, the bandwidth remained fairly constant for switch as well as client despite varying offloading parameters when flows traversed through an OpenFlow switch.

In figavgcpuclientnoswitch, x-axis represents the offloading parameters, while y-axis represents the average percent CPU utilization for TCP and UDP flows for a client when flows directly traverse from client to server. Irrespective of offloading settings, 100% CPU utilization is observed for UDP flows. This makes sense, since no offloading should occur for UDP flows. However, we see almost 100% CPU utilization for TCP flows in the absence of TSO/GSO. This is in line with findings expressed in figavgBandClientNoswitch. Since the NIC non longer performs offloading, CPU is utilized for segmentation. The CPU utilization increases 10x with a 2x decrease in bandwidth. In figavgcpuclientswitch, x-axis represents the offloading parameters, while y-axis represents the average percent CPU utilization for TCP and UDP flows when a switch is placed in server-client path. As observed in figavgcpuclientnoswitch, the average CPU utilization remains unaffected by varying offloading parameters for UDP flows. However, for TCP flows, a significantly uniform and small percent of CPU utilization is observed.

In figavgcpuservernoswitch, the x-axis represents the offloading parameters, while y-axis represents the average percent cpu utilization. The fairly even distribution of cpu percent utilization with or without switch suggests that packet segmentation does not occur at the server end.

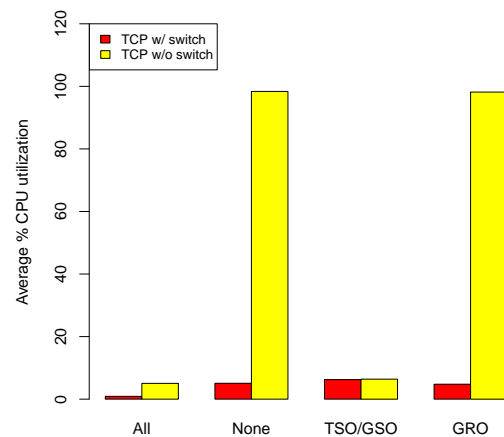


Figure 1: The x-axis represents offloading parameters while y-axis looks at average percent CPU utilization

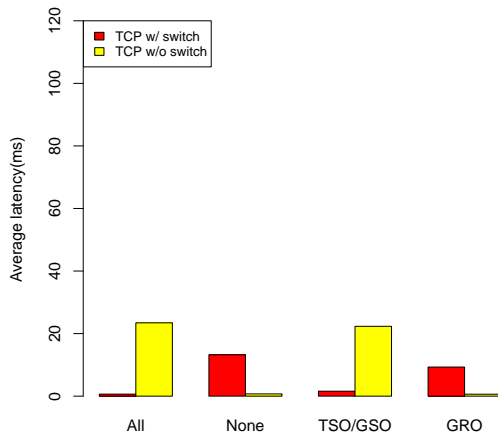


Figure 2: The x-axis represents offloading parameters while y-axis looks at average latency

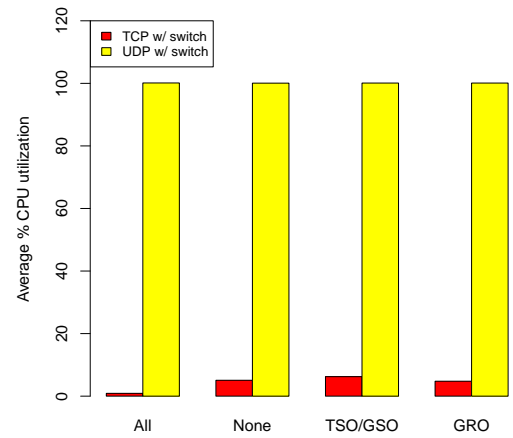


Figure 4: The x-axis represents offloading parameters while y-axis looks at average percent CPU utilization in presence of switch on client end

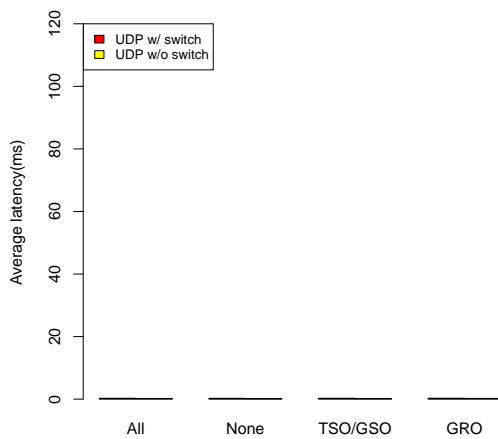


Figure 3: The x-axis represents offloading parameters while y-axis looks at average latency

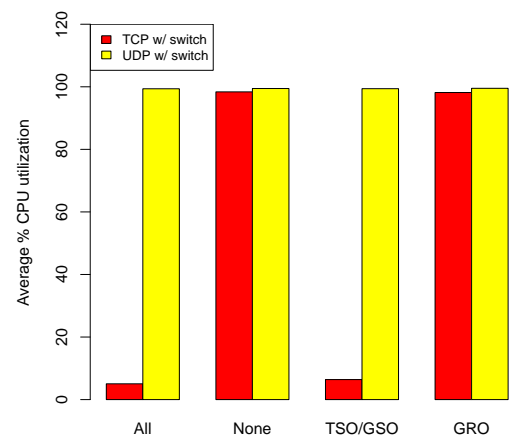


Figure 5: The x-axis represents offloading parameters while y-axis looks at average percent CPU utilization in absence of switch on client end

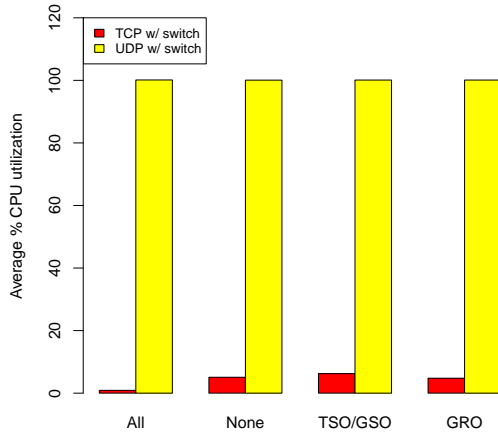


Figure 6: The x-axis represents offloading parameters while y-axis looks at average percent CPU utilization in presence of switch on server end

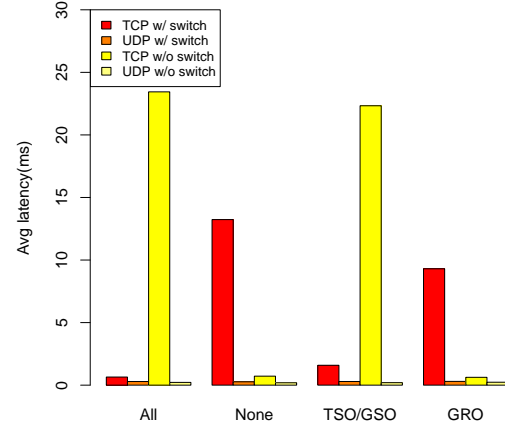


Figure 8: The x-axis represents offloading parameters while y-axis looks at average latency for FIFO algorithm

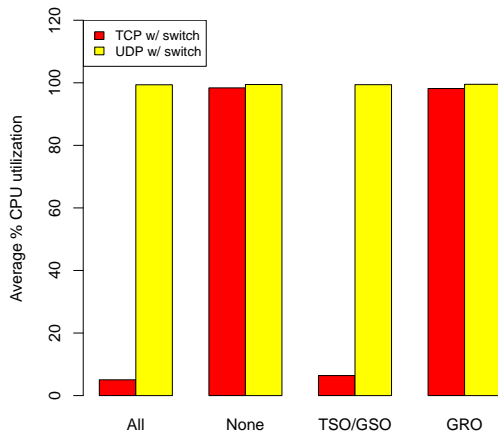


Figure 7: The x-axis represents offloading parameters while y-axis looks at average percent CPU utilization in absence of switch on server end

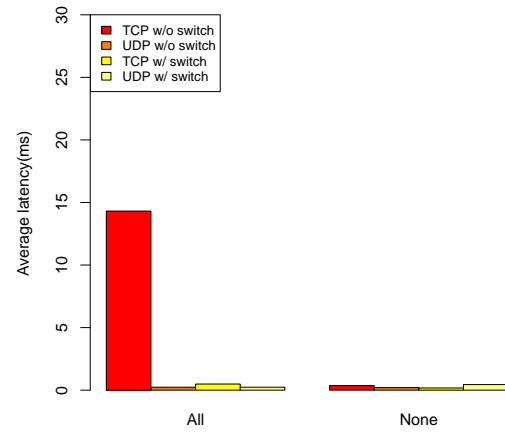


Figure 9: The x-axis represents offloading parameters while y-axis looks at average latency for SFQ algorithm

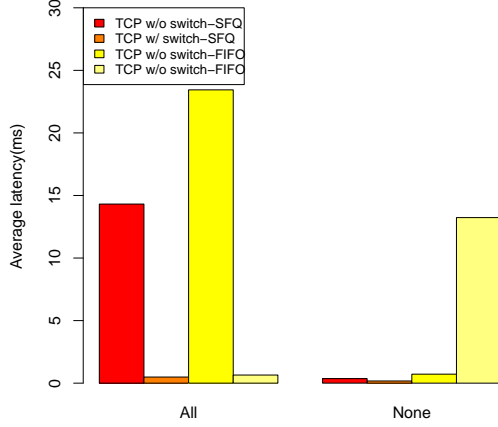


Figure 10: The x-axis represents offloading parameters while y-axis looks at average latency for SFQ algorithm with TCP flows

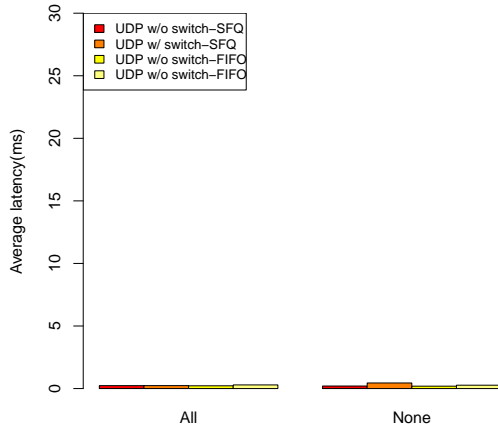


Figure 11: The x-axis represents offloading parameters while y-axis looks at average latency for SFQ algorithm with UDP flows

3.2 Latency

In `figavglatencyfifosfq`, the x-axis represents the offloading parameters, while y-axis represents the latency in ms for TCP flows at the client end. A distinct difference in latency is observed for when GRO is switched off for SFQ. In `figavglatencyFIFO`, we see the average latency in ms for UDP flows is independent of offloading settings.

However, for TCP, TSO/GSO's presence dramatically increases latency in absence of switch.

While studying the latency profile for UDP flows, we notice a significantly higher latency in presence of TSO/GSO for without a switch case. However, not sure how this makes sense.

4 Limitations and Future Work

Remote Direct Memory Access (RDMA) NICs would be interesting to study since it simplifies many problems observed with generic TCP offloading. It would be especially interesting and important to relate latency observed from utilizing RDMA and CPU utilization on non-RDMA based NICs. Since RDMA enables the network adapter to transfer data directly to or from application memory, no work is required from CPUs, caches or context switches.

Maximum Transmission Unit (MTU) is the greatest amount of data that can be transferred in one physical frame over the network. Thus, increasing MTU, may decrease latency. The standard MTU for ethernet is 1500 bytes. We would have like to increase that upto 9000 units, but were unable to perform experiments due to hardware limitations.

5 Conclusion

6 Acknowledgments

We would like to thank Aaron Gember without whose help, running experiments on the Open Flow test best at the University of Wisconsin-Madison would not be possible.

References

- [1] Dctcp linux kernel patch.
- [2] ALIZADEH, M., KABBANI, A., EDSALL, T., PRABHAKAR, B., VAHDAT, A., AND YASUDA, M. Less is more: trading a little bandwidth for ultra-low latency in the data center. In *Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation* (Berkeley, CA, USA, 2012), NSDI'12, USENIX Association, pp. 19–19.
- [3] CHIOLA, G., AND CIACCIO, G. Architectural issues and preliminary benchmarking of a low-cost network of workstations based on active messages. In *In Proc. 14th ITG/GI conference on Architecture of Computer Systems (ARCS'97)* (1997), pp. 13–22.
- [4] CHUN, B. N., MAINWARING, A. M., AND CULLER, D. E. Virtual network transport protocols for myrinet. In *IEEE MICRO* (1997), pp. 53–63.
- [5] DEMERS, A., KESHAV, S., AND SHENKER, S. Analysis and simulation of a fair queueing algorithm. *SIGCOMM Comput. Commun. Rev.* 19, 4 (Aug. 1989), 1–12.
- [6] DITTIA, Z. D. Integrated hardware/software design of a high-performance network interface, 2001.

- [7] D.LACAMERA. Tcp pacing linux implementation.
- [8] EICKEN, T. V., BASU, A., BUCH, V., AND VOGELS, W. U-net: A user-level network interface for parallel and distributed computing. In *In Fifteenth ACM Symposium on Operating System Principles* (1995).
- [9] MCKEOWN, N., ANDERSON, T., BALAKRISHNAN, H., PARULKAR, G., PETERSON, L., REXFORD, J., SHENKER, S., AND TURNER, J. Openflow: enabling innovation in campus networks. *SIGCOMM Comput. Commun. Rev.* 38, 2 (Mar. 2008), 69–74.