

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 present a measurement study where we examine various factors such as CPU utilization, send and receive offloading, and how queueing affects latency.

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 up to 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 examining ways of reducing latency in data centers is interesting as the data center operator has control of all aspects of the network. However, it is still important to investigate end-host solutions for cases where a tenant is co-locating servers in a data center and does not control the network topology or switches. Reducing latency is especially important in 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 by never allowing queues to build up at the switch. In this paper, we aim to understand how offloading, switches, NICs, and queues affect latency, and what are the potential trade-offs.

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. We also aim to examine the effects of queueing. In one set of experiments, we use TCP which by nature will cause queueing. In the other set of experiments, we use a 500 Mbps UDP flow to

simulate some traffic on the network, but latency sensitive flows should not be queued. To represent a latency sensitive flow, we ping the other host every 200 ms for the duration of the TCP/UDP flow. 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 outbound 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. We conduct our experiments with the following offloading configurations at end-hosts:

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 (sender of traffic) and the other as server (receiver of traffic). Eight experiments were run in this setting. Another eight experiments were run with the two machines connected directly together. For each experiment, throughput is recorded using *iperf*, CPU utilization through *pidstat*, inter-packet delay with *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. These experiments were run on multi-core machines at coarse timescales (1 second for CPU measurements) and it is our belief that our instrumentation adds little overhead.

3 Evaluation

3.1 Bandwidth

TSO works by moving the work of breaking down large chunks of data into MTU-sized packets over to the NIC. In *figavgBandwidthClientNoswitch*, the x-axis represents offloading settings while y-axis represents bandwidth. Bandwidth is almost halved when TSO/GSO is

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 no 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.

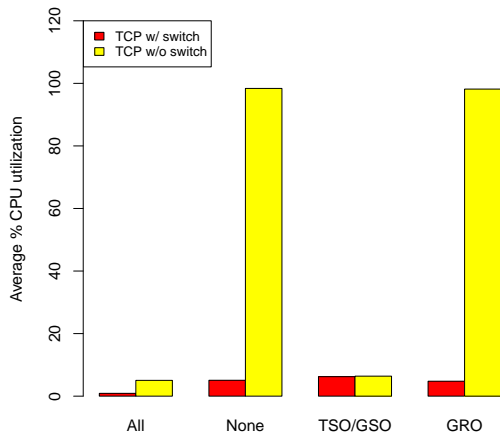


Figure 1: The x-axis represents offloading parameters while y-axis looks at average percent CPU utilization

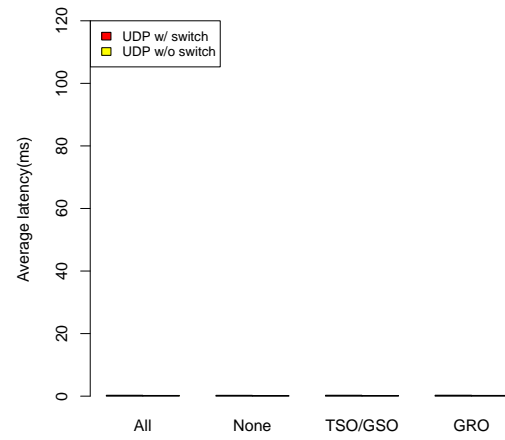


Figure 3: The x-axis represents offloading parameters while y-axis looks at average latency

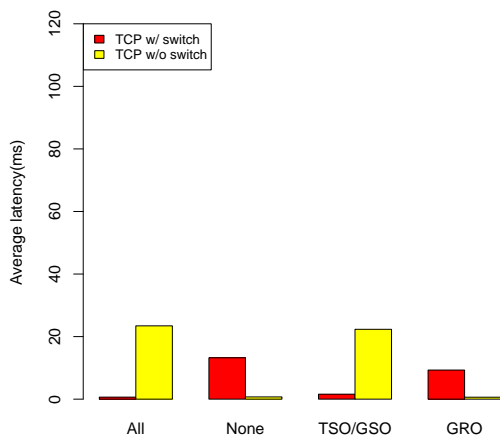


Figure 2: The x-axis represents offloading parameters while y-axis looks at average latency

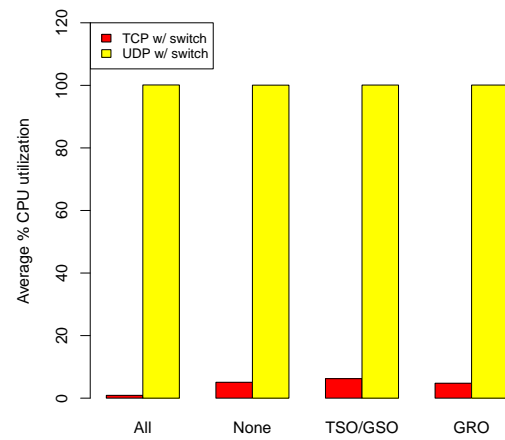


Figure 4: The x-axis represents offloading parameters while y-axis looks at average percent CPU utilization in presence of switch on client end

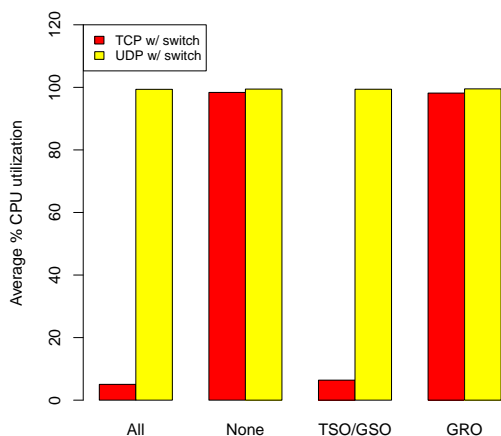


Figure 5: The x-axis represents offloading parameters while y-axis looks at average percent CPU utilization in absence of switch on client end

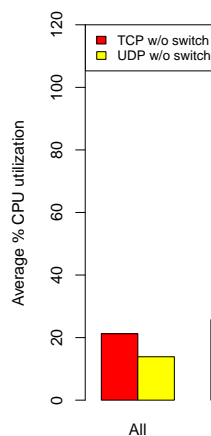


Figure 7: The x-axis represents offloading parameters while y-axis looks at average percent CPU utilization in absence of switch on server end

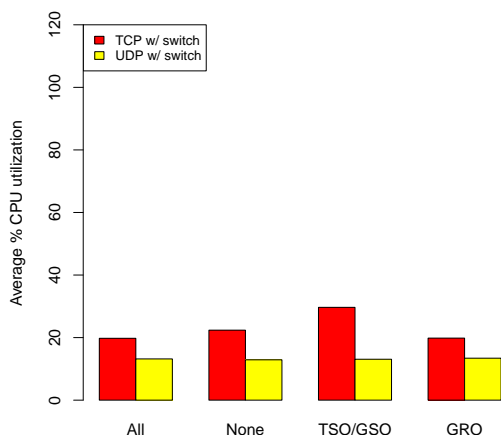


Figure 6: The x-axis represents offloading parameters while y-axis looks at average percent CPU utilization in presence of switch on server end

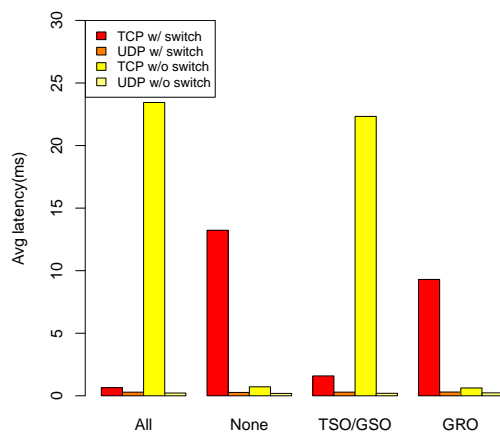


Figure 8: The x-axis represents offloading parameters while y-axis looks at average latency for FIFO algorithm

These results in Figure 8 would indicate that there were larger buffers at end-hosts than the switch. While there aren't any data center topologies that connect end-hosts directly together, this illustrates the importance of sizing the buffers correctly. Having all offloading reduces latency. Avg UDP latency w/ switch is .29-.3ms. UDP latency w/o switch is .191-.234ms. Roughly .1ms for a RTT.

TSO works by moving work of breaking down large

chunks of data over to the NIC. In Figure 12, 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.

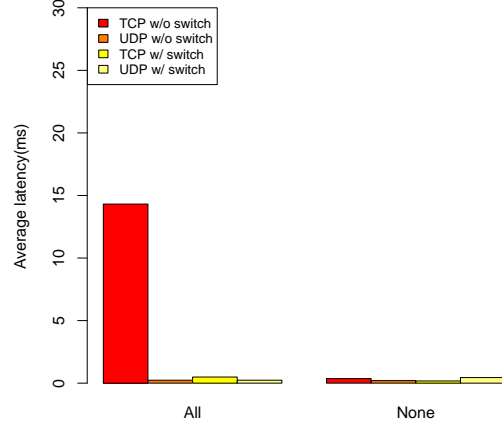


Figure 9: The x-axis represents offloading parameters while y-axis looks at average latency for SFQ algorithm

In Figure ??, 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 Figure 12. Since the NIC no longer performs offloading, CPU is utilized for segmentation. The CPU utilization increases 10x with a 2x decrease in bandwidth. In ??, 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 ??, 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.

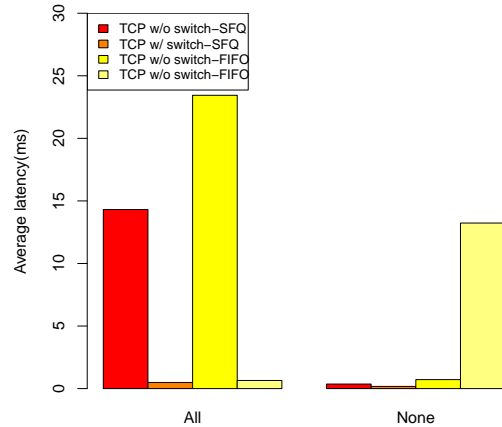


Figure 10: The x-axis represents offloading parameters while y-axis looks at average latency for SFQ algorithm with TCP flows

In ??, 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.

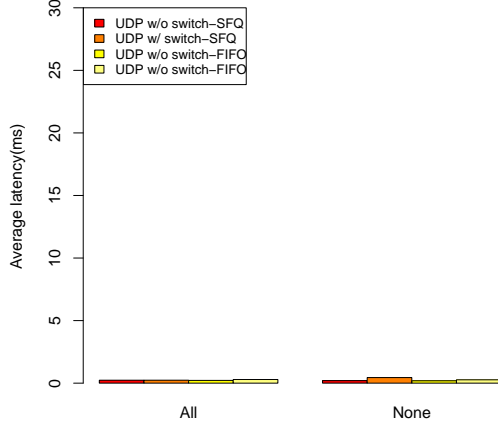


Figure 11: The x-axis represents offloading parameters while y-axis looks at average latency for SFQ algorithm with UDP flows

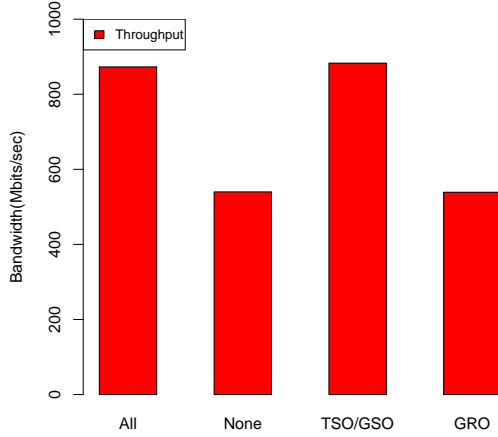


Figure 12: The x-axis represents offloading parameters while y-axis looks at average bandwidth in mbits/sec on client end in absence of switch

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 liked to increase that up to 9000 units, but were unable to perform experiments due to hardware limitations.

5 Conclusion

In this paper we study types of flows, influence of switches and varying offloading parameters on TCP. We show that optimizing for buffer sizes is important to achieve desired goals. Offloading helps reduce latency in TCP flows, while UDP flows remain largely uninfluenced. Throughput remains mostly unchanged for TCP, except in the no switch case. Throughput dropped when TSO/GSO is disabled. This is likely due to 100% utilization, but it's unclear in what's different between the switch/no switch case. SFQ provides lower average latency, but higher variation. 15ms stdev when all settings are enabled. Increases in CPU utilization seem minimal/non-existent. These

6 Acknowledgments

We would like to thank Aaron Gember without whose help, running experiments on the Open Flow test bed 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.