

High Performance Memcached

Lijing Chen

4th Year Project Report
Computer Science
School of Informatics
University of Edinburgh

2018

Abstract

This is an example of `infthesis` style. The file `skeleton.tex` generates this document and can be used to get a “skeleton” for your thesis. The abstract should summarise your report and fit in the space on the first page. You may, of course, use any other software to write your report, as long as you follow the same style. That means: producing a title page as given here, and including a table of contents and bibliography.

Acknowledgements

Acknowledgements go here.

Table of Contents

1	Introduction	7
1.1	Motivation	7
2	Background	9
2.1	Memcached	9
2.1.1	Memcached protocol	9
2.2	User-level networking	9
2.2.1	Data Plane Development Kit	10
2.2.2	Seastar	10
2.3	TCP Offload	10
2.3.1	Full offload vs Partial offload	10
2.3.2	TCP/IP checksum offload	11
2.3.3	Large send offload	11
3	Methodology	13
3.1	Quality of Service	13
3.2	Hardware	13
3.2.1	Server	13
3.3	Benchmark	14
3.3.1	Open-loop vs Closed loop	14
3.3.2	Treadmill	14
3.3.3	Memaslap	14
4	Baseline Tuning	15
4.1	Multithreading	15
4.2	Multiprocessing	15
4.3	IRQ Affinity	15
5	Contention Analysis	17
5.1	Core contention	17
5.2	SMT Thread contention	17
6	User-level Networking Optimization	19
6.1	DPDK	19
6.1.1	DPDK Setup	19
7	Evaluation	21

8	TCP Offloading	23
8.1	Setup TCP Offloading	23
8.2	Large Send Offload	23
9	Performance Analysis	25
9.1	Profiling Tools	25
9.1.1	Linux perf	25
9.2	Profiling Visualization	25
9.2.1	FlameGraph	25
	Bibliography	27

Chapter 1

Introduction

1.1 Motivation

Software object caches, such as Memcached, are crucial to achieving high throughput and low response latency in a datacenter setting. A well-known problem with Memcached deployments is the high overhead incurred by the kernel's TCP/IP stack in processing of individual requests. This project will evaluate the performance of a vanilla Memcached setup, followed by an assessment of the effect that user-level networking and TCP offload have on Memcached.

Memory object caches is now widely used in industry. One of the most widely used memory object cache is Memcached.

Chapter 2

Background

2.1 Memcached

Memcached is a "high-performance, distributed memory object caching system, generic in nature, but intended for use in speeding up dynamic web applications by alleviating database load." Despite the official description aimed at dynamic web applications, Memcached is also used as a generic key value store to locate servers and services .

2.1.1 Memcached protocol

Memcached support two transport-layer protocols: TCP and UDP. They have some small difference. The basic protocol we are going to use in our experiments are described as Table 2.1. Note that there are many other command provided by memcached, but we are only going to test the most basic ones as listed. Comparing with TCP, the UDP protocol added an additional binary header for controlling.

2.2 User-level networking

One performance bottleneck of modern high-performance applications is the kernel TCP/IP stack. There is a trend in industry to replace kernel TCP/IP stack with user-level TCP/IP stack.

Operation	Meaning
get [key]	Get the value of key
set [arg1] [flags] [xxx]	Mellanox Connect-X 4
delete [key]	delete a key provided by the argument.

Table 2.1: Basic Memcached protocol

2.2.1 Data Plane Development Kit

The Data Plane Development Kit is "a set of libraries and drivers for fast packet processing"[3]. It runs mostly in userspace, and it serves as a platform for developer to develop high performance networking applications.

DPDK provides developers with a framework, which contains a set of libraries for different hardware and software environments. DPDK encapsulates the low level environmental detail by creating Environment Abstraction Layer (EAL). Developers can program with a general interface, and link the library compiled for each specific environment, instead of coding with specific API from different hardware, operating systems etc.

2.2.2 Seastar

DPDK doesn't contain a TCP/IP stack. In order to migrate TCP/IP based applications to use DPDK, developers need to implement their own TCP/IP stack using DPDK provided API.

2.3 TCP Offload

TCP Offloading is another technique quite widely used in modern NIC cards. In order to alleviate the load of CPU, NIC manufacturers provided different level of offloading features in NIC, this means some of the operations like checksum, fragmentation can be processed by NIC directly instead of host CPU.

2.3.1 Full offload vs Partial offload

There are two type of TCP Offloading. Full offload means we offload all kernel TCP/IP stack to our NIC. Partial offload means we only offload some part of TCP/IP stack to our NIC, like checksum, segmentation etc.

TOE implementations also can be differentiated by the amount of processing that is offloaded to the network adapter. In situations where TCP connections are stable and packet drops infrequent, the highest amount of TCP processing is spent in data transmission and reception. Offloading just the processing related to transmission and reception is referred to as partial offloading. A partial, or data path, TOE implementation eliminates the host CPU overhead created by transmission and reception. However, the partial offloading method improves performance only in situations where TCP connections are created and held for a long time and errors and lost packets are infrequent. Partial offloading relies on the host stack to handle control that is, connection setups as well as exceptions. A partial TOE implementation does not handle the following:

2.3.2 TCP/IP checksum offload

TCP/IP checksum offload let the network adapter to do the calculation task of verifying the checksum of packet received, which is done by CPU originally. This technique can reduce the CPU utilization.

The TCP/IP checksum offload technique moves the calculation of the TCP and IP checksum packets from the host CPU to the network adapter. For the TCP checksum, the transport layer on the host calculates the TCP pseudo-header checksum and places this value in the checksum field, thus enabling the network adapter to calculate the correct TCP checksum without touching the IP header. However, this approach yields only a modest reduction in CPU utilization.

2.3.3 Large send offload

Large send offload (LSO), also known as TCP segmentation offload (TSO), frees the OS from the task of segmenting the applications transmit data into MTU-size chunks. Using LSO, TCP can transmit a chunk of data larger than the MTU size to the network adapter. The adapter driver then divides the data into MTU-size chunks and uses the prototype TCP and IP headers of the send buffer to create TCP/IP headers for each packet in preparation for transmission.

LSO is an extremely useful technology to scale performance across multiple Gigabit Ethernet links, although it does so under certain conditions. The LSO technique is most efficient when transferring large messages. Also, because LSO is a stateless offload, it yields performance benefits only for traffic being sent; it offers no improvements for traffic being received. Although LSO can reduce CPU utilization by approximately half, this benefit can be realized only if the receivers TCP window size is set to 64 KB. LSO has little effect on interrupt processing because it is a transmit-only offload.

Methods such as TCP/IP checksum offload and LSO provide limited performance gains or are advantageous only under certain conditions. For example, LSO is less effective when transmitting several smaller-sized packages. Also, in environments where packets are frequently dropped and connections lost, connection setup and maintenance consume a significant proportion of the hosts processing power. Methods like LSO would produce minimal performance improvements in such environments.

Chapter 3

Methodology

3.1 Quality of Service

Firstly, it is important to define the desired quality of service we are looking to target with our benchmarks. Frequently, distributed systems are designed to work in parallel, each component responsible for a piece of computation which is then ultimately assembled into a larger piece of response before being shipped to the client. For example, an e-commerce store may choose to compute suggested products as well as brand new products separately only to assemble individual responses into an HTML page. Therefore, the slowest of all individual components will determine the overall time required to render a response.

Let us define the quality of service (QoS) target of this study. For our benchmarking purposes, a sufficient QoS will be the 99th percentile tail latency of a system under 1 millisecond. This is a reasonable target as the mean latency will generally (based on latency distribution) be significantly smaller. Furthermore, it is a similar latency target used in related research.

3.2 Hardware

The Server we are using is described in Table 3.1.

3.2.1 Server

The Network Interface Card we are using is (mlx5) [1] .

Hardware List	
Component Name	Device
CPU	Intel(R) Xeon(R) CPU E5-2630 v4 @ 2.20GHz (40 cores)
Network Interface Card	Mellanox Connect-X 4

Table 3.1: Hardware Table

3.2.1.1 Network topology

For DPDK tests, we directly connect the NIC in two servers using network cable due to device limitation. The reason we use this topology is the NIC only supports Infiniband ports, but we don't have a switch that supports Infiniband. But then don't make a large difference and it's sufficient for us to do our experiments.

3.3 Benchmark

3.3.1 Open-loop vs Closed loop

There are two kinds of load testers: open-loop and closed-loop. Open-loop load testers send a new request only when the previous one has returned (only one outstanding request), while closed-loop load testers can send a new one even the previous one is pending (allow multiple outstanding requests).

3.3.2 Treadmill

Treadmill [4] is a open-source load tester for memcached developed by Facebook. It's a typical open-loop load tester. It only supports TCP for memcached. For our experiments involving UDP, we implemented our own version of memcached for UDP.

3.3.3 Memaslap

Memaslap is a load tester provide by libmemcached [2]. It's a closed loop load tester.

Chapter 4

Baseline Tuning

4.1 Multithreading

In this section, we are going to explore the influence of multithreading.

4.2 Multiprocessing

4.3 IRQ Affinity

Chapter 5

Contention Analysis

5.1 Core contention

5.2 SMT Thread contention

Chapter 6

User-level Networking Optimization

6.1 DPDK

6.1.1 DPDK Setup

In this section, we are going to go over how DPDK is set up.

Setup hugepage mapping Hugepages is a mechanism that allows the Linux kernel to utilize the multiple page size capabilities of modern hardware architectures. Linux uses pages as the basic unit of memory, where physical memory is partitioned and accessed using the basic page unit. The default page size is 4096 Bytes in the x86 architecture. Hugepages allows large amounts of memory to be utilized with a reduced overhead.

Compile and insert DPDK kernel modules First we need to do `make config` to tell the compiler which architecture we are targeting. We need to generate binaries working on `x86_64-native-linuxapp-gcc`.

```
$ make config T=x86_64-native-linuxapp-gcc
```

Some NIC driver is not enabled by default, including the one we are using (mlx5). In this case, we need to manually enable it in `config/common_base` file.

Then we can use `make` command to compile

Then switch to output folder and insert kernel module: `igb_uio.ko` using `insmod`.

Configure NIC First we need to unload the NIC from system kernel driver using `ifconfig`. `ifconfig p2p1 down`. Then we can use the tool provided by DPDK to bind NIC to kernel driver. `dpdk-devbind.py`

Chapter 7

Evaluation

Chapter 8

TCP Offloading

8.1 Setup TCP Offloading

First, we use **ethtool** to check device offloading. We can show the support for TCP Offloading by executing following command:

```
# ethtool -k p2p2
```

8.2 Large Send Offload

Chapter 9

Performance Analysis

9.1 Profiling Tools

9.1.1 Linux perf

9.2 Profiling Visualization

9.2.1 FlameGraph

Bibliography

- [1] Hiroki Arimura. Learning acyclic first-order horn sentences from entailment. In *Proc. of the 8th Intl. Conf. on Algorithmic Learning Theory, ALT '97*, pages 432–445, 1997.
- [2] Chen-Chung Chang and H. Jerome Keisler. *Model Theory*. North-Holland, third edition, 1990.
- [3] Intel. Intel Data Plane Development Kit. <http://www.intel.com/go/dpdk>, 2016.
- [4] Yunqi Zhang, David Meisner, Jason Mars, and Lingjia Tang. Treadmill: Attributing the source of tail latency through precise load testing and statistical inference. In *43rd ACM/IEEE Annual International Symposium on Computer Architecture, ISCA 2016, Seoul, South Korea, June 18-22, 2016*, pages 456–468, 2016.