

Memcached vs Redis: Benchmarking In-memory Object Caches

Milan Pavlik

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

2016

Abstract

0.1 Abstract

Abstract Goes here

Acknowledgements

I would like to express my gratitude to my supervisor, Dr. Boris Grot, for his guidance, patience and experience provided during my dissertation.

I would like thank my parents, sister and friends for their continuous support, encouragement and willingness to listen.

Table of Contents

0.1	Abstract	3
1	Introduction	7
1.1	Motivation	7
1.2	Memory Object Caches	7
1.2.1	Desired qualities	7
1.2.2	Design and Implementations	8
1.2.3	Performance metrics	8
1.2.4	Memcached	9
1.3	Methodology	10
1.3.1	Quality of Service	11
1.3.2	Hardware	11
1.3.3	Workload generation	11
1.3.4	Benchmark	12
2	Memcached	15
2.1	Default Performance	15
2.1.1	Default Throughput vs Latency	16
2.1.2	Effect of Connections	16
2.1.3	Server CPU	17
2.2	Memcached Thread Scalability	18
2.2.1	Throughput & Latency	19
2.2.2	CPU Time	20
2.2.3	Thread evaluation	20
2.3	Thread pinning	21
2.3.1	Throughput vs Latency	21
2.4	Group Size	22
2.5	Receive & Transmit Queues fixing	23
2.6	Multiple Memcached Processes	23
3	Redis	25
3.1	Out of the Box Performance	25
3.1.1	Latency vs Throughput	26
3.1.2	CPU Utilization	26
3.2	Multiple Redis Instances	27
3.2.1	Latency and Throughput	28
3.2.2	CPU Utilization	30

3.3	Pinned Redis Instances	31
3.3.1	Pinned Latency and Throughput	31
3.3.2	Redis Persistence	32
3.4	Object Size	32
3.4.1	Latency and Throughput	32
3.4.2	CPU Utilization	34
3.5	Key Distributions	34
3.5.1	Gaussian distribution	35
3.5.2	Zipf distribution	35
4	Redis & Memcached: Head to Tail	37
5	Conclusion	39
	Bibliography	41

Chapter 1

Introduction

1.1 Motivation

TODO

1.2 Memory Object Caches

Firstly, the purpose of a memory object cache is to use the machine's available RAM for key-value storage. The implication of a *memory object cache* is that data is only stored in memory and should not be offloaded on the hardware in order to not incur hard drive retrieval delay. As a result, memory caches are often explicitly configured with the maximum amount of memory available.

Secondly, an *object* cache implies that the cache itself is not concerned with the type of data (binary, text) stored within. As a result, memory object caches are multi-purpose caches capable of storage of any data type within size restrictions imposed by the cache.

Finally, memory object caches can be deployed as single purpose servers or also co-located with another deployment. Consequently, general purpose object caches often provide multiple protocols for accessing the cache - socket communication or TCP over the network. Both caches in question - Memcached and Redis - support both deployment strategies. Our primary focus will be on networked protocols used to access the cache.

1.2.1 Desired qualities

Firstly, an object cache should support a simple interface providing the following operations - *get*, *set* and *delete* to retrieve, store and invalidate an entry respectively.

Secondly, a general purpose object cache should have the capability to store items of arbitrary format and size provided the size satisfies the upper bound size constraints imposed by the cache. Making no distinction between the type of data is a fundamental generalization of an object cache and allows a greater degree of interoperability.

Thirdly, a cache should support operation atomicity in order to prevent data corruption resulting from multiple simultaneous writes.

Furthermore, cache operations should be performed efficiently, ideally in constant time and the cache should be capable of enforcing a consistent eviction policy in the case of memory bounds are exceeded.

Finally, a general purpose object cache should be capable of handling a large number of requests per second while maintaining a fair and as low as possible quality of service for all connected clients.

1.2.2 Design and Implementations

The design and implementation of a general purpose cache system is heavily influenced by the desired qualities of a cache.

Firstly, high performance requirement and the need for storage of entries of varying size generally requires the cache system to implement custom memory management models. As a result, a mapping data structure with key hashing is used to efficiently locate entries in the cache.

Secondly, due to memory restrictions, the cache is responsible for enforcing an eviction policy. Most state of the art caches utilize least recently used (LRU) cache eviction policy, however, other policies such as first-in-first-out can also be used.

In the case of *Memcached*, multi-threaded approach is utilized in order to improve performance. Conversely to *Memcached*, *Redis* is implemented as a single threaded application and focuses primarily on a fast execution loop rather than parallel computation.

1.2.3 Performance metrics

Firstly, the primary metrics reflecting performance of an in memory object cache are *mean latency*, *99th percentile latency* and *throughput*. Both latency statistics are reflective of the quality of service the cache is delivering to it's clients. Throughput is indicative of the overall load the cache is capable of supporting, however, throughput is tightly related to latency and on it's own is not indicative of the real cache performance under quality constraints.

Secondly, being a high performance application with potentially network, understanding the proportion of CPU time spent inside the cache application compared to time

spent processing network requests and handling operating system calls becomes important. Having an insight into the CPU time breakdown allows us to better understand bottlenecks of the application.

Finally, the *hit* and *miss* rate of the cache can be used as a metric, particularly when evaluating a cache eviction policy, however, the hit and miss rate is tightly correlated with the type of application and the application context and therefore it is not a suitable metric for evaluating performance alone.

1.2.4 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.” [4] Despite the official description aimed at dynamic web applications, memcached is also used as a generic key value store to locate servers and services [1].

1.2.4.1 Memcached API

Memcached provides a simple communication protocol. It implements the following core operations:

- `get key1 [key2..N]` - Retrieve one or more values for given keys,
- `set key value [flag] [expiration] [size]` - Insert *key* into the cache with a *value*. Overwrites current item.
- `delete key` - Delete a given key.

Memcached further implements additional useful operations such as `incr/decr` which increments or decrements a value and `append/prepend` which append or prepend a given key.

1.2.4.2 Implementation

Firstly, Memcached is implemented as a multi-threaded application. “Memcache instance started with *n* threads will spawn *n* + 1 threads of which the first *n* are worker threads and the last is a maintenance thread used for hash table expansion under high load factor.” [3]

Secondly, in order to provide performance as well as portability, memcached is implemented on top of *libevent* [9]. “The libevent API provides a mechanism to execute a callback function when a specific event occurs on a file descriptor or after a timeout has been reached. Furthermore, libevent also support callbacks due to signals or regular timeouts.” [9]

Thirdly, Memcached provides guarantees on the order of actions performed. Therefore, consecutive writes of the same key will result in the last incoming request being the retained by memcached. Consequently, all actions performed are internally atomic.

As a result, memcached employs a locking mechanism in order to be able to guarantee order of writes as well as execute concurrently. Internally, the process of handling a request is as follows:

1. Requests are received by the Network Interface Controller (NIC) and queued
2. *Libevent* receives the request and delivers it to the memcached application
3. A worker thread receives a request, parses it and determines the command required
4. The *key* in the request is used to calculate a hash value to access the memory location in $O(1)$
5. Cache lock is acquired (*entering critical section*)
6. Command is processed and LRU policy is enforced
7. Cache lock is released (*leaving critical section*)
8. Response is constructed and transmitted [12]

We can observe that steps 1-4 and 8 can be parallelized without the need for resource locking. However, the critical section in steps 5-7 is executed with the acquisition of a global lock. Therefore, at this stage execute is not being performed in parallel.

1.2.4.3 Production deployments

TODO: Discuss Facebook, Amazon, Twitter, ... deployments of memcached

1.3 Methodology

In order to effectively benchmark the performance of both types of caches in question, it is essential to be able to stress the cache server sufficiently to experience queuing delay and saturate the server. This study is concerned with the performance of the cache server rather than performance of the underlying network and therefore it is essential to utilize a sufficient number of clients in order to saturate the server while maintaining low congestion on the underlying network.

The benchmarking methodology is heavily influenced by similar studies and benchmarks in the literature. This allows for a comparison of observed results and allows for a better correlation with related research.

1.3.1 Quality of Service

Firstly, it is important the desired quality of service we are looking to benchmark for. 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 [6].

1.3.2 Hardware

Performance benchmarks executed in this study will be run on 8 distinct machines with the following configuration: *6 core Intel(R) Xeon(R) CPU E5-2603 v3 @ 1.60GHz*, *8 GB RAM* and *1Gb/s Network Interface Controller (NIC)*.

All the hosts are connected to a *Pica8 P-3297* switch with 48 1Gbps ports arranged in a star topology. A single host is used to run an object cache system while the remaining seven are used to generate workloads against the server.

1.3.3 Workload generation

Workload for the cache server is generated using Memtier Benchmark developed by Redis Labs [5]. Memtier has been chosen as the benchmark for this study due to its high level of configurability as well as ability to benchmark both *Memcached* and *Redis*. Utilizing the same benchmark client for both caches allows for a decreased variability in results when a comparison is made.

In order to create a more realistic simulation of a given workload, 7 servers all running *memtier* simultaneously are used. A simple parallel ssh utility is used to start, stop and collect statistics from the load generating clients.

The workload generated by Memtier is driven by the configuration specified. The keys and values are drawn from a configured distribution dynamically at runtime. All comparable benchmarks presented in this thesis configure the same initial seed for comparable benchmarks in order to minimize stochastic behavior.

1.3.4 Benchmark

In the context of this thesis, a benchmark is a set of workloads executed against the cache host. Statistics are collected from the cache host as well as the clients in order to draw conclusions.

Firstly, a benchmark consists of a warm up stage. The cache is being loaded with initial data in order to prevent a large number of cache misses and skewed results.

Secondly, a configured workload is generated and issued against the cache host from multiple benchmarking hosts simultaneously.

Thirdly, the workload is repeated 2 more times in order to decrease the impact of stochastic events in the benchmark.

Finally, statistics are collected, individual benchmark runs are averaged and the results are processed.

1.3.4.1 Memtier

Memtier benchmark is “a command line utility developed by Redis Labs for load generation and benchmarking NoSQL key-value databases” [5]. It provides a high level of configurability allowing for example to specify patterns of *sets* and *gets* as well as generation of key-value pairs according to various distributions, including Gaussian and pseudo-random.

Memtier is a threaded application built on top of `libevent` [9], allowing the user to configure the number of threads as well as the number of connections per each thread which can be used to control the server load. Additionally, memtier collects benchmark statistics including latency distribution, throughput and mean latency. The statistics reported are used to draw conclusions on the performance under a given load.

Memtier execution model is based on the number of threads and connections configured. For each thread t , there are c connections created. The execution pattern within each thread is as follows:

1. Initiate c connections
2. For each connection
 - (a) Make a request to the cache server
 - (b) Provide a *libevent* callback to handle response outside of the main event loop
3. Tear down c connections

By offloading response handling to a callback inside `libevent`, memtier is able to process a large number of requests without blocking the main event loop until a response from the network request is returned while maintaining the ability to collect statistics effectively.

Connections created with the target server are only destroyed at the end of the benchmark. This is a realistic scenario as in a large distributed environment the cache clients will maintain open connections to the cache to reduce the overhead of establishing a connection.

Mentier provides a comprehensive set of configuration options to customize Mentier behavior and tailor the load. Table 1.1 outlines the relevant configuration options. The complete set of configuration options is available on RedisLabs [11].

Configuration option	Explanation	Default Value
-s	Server Address	localhost
-p	Port number	6379
-P	Protocol - redis, memcache_text, memcache_binary	redis
-c	Number of Clients per Thread	50
-t	Number of Threads	4
-data-size	The size of the object to send in bytes	32
-random-data	Data should be randomized	false
-key-minimum	The minimum value of keys to generate	0
-key-maximum	The maximum value of keys to generate	10 million

Table 1.1: Mentier Configuration Options

1.3.4.2 Open-loop vs Closed-loop

A load tester can be constructed with different architecture in mind. The main two types of load testers are *open-loop* and *closed-loop*. Closed-loop load testers frequently construct and send a new request only when the previous request has received a response. On the other hand, open-loop principle aims to send requests in timed intervals regardless of the response from the previous requests.

The consequence of a closed-loop load tester is potentially reduced queuing on the server side and therefore observed latency distribution may be lower than when server side queuing is observed.

Mentier falls in the category of closed loop testers when considering a single thread of mentier. However, mentier threads are independent of each other and therefore requests for another connection are made even if the previous request has not responded. Furthermore, by running mentier on multiple hosts simultaneously, the closed loop implications are alleviated and the server observes queuing delay in the network stack.

Chapter 2

Memcached

The purpose of this chapter is to benchmark and evaluate memcached performance. Firstly, we will examine performance under default configuration of both the server and the client. Secondly, threading will be explored in relation to latency and throughput. Thirdly, the effect of memcached's `group size` will be explored in relation to performance. Additionally, configuration of receive and transmit queues will be explored and finally, an execution model of multiple processes will be visited in order to establish a comparison baseline. Throughout the benchmarks, we will be focusing cache performance which meets desired the QoS.

2.1 Default Performance

Firstly, it is essential to establish a performance baseline of *memcached* under high utilization. In order to establish the baseline, a default configuration of memcached will be used with the exception of the amount of memory allocated for exclusive use by the application. The *memcached* application will be started with

```
memcached -d -p 11120 -m 6144
```

specifying the port and the amount of memory to be used by the application.

To find a saturation point, we can increase the number of connections linearly and analyze the results. To load test the cache, we execute the following command on all client servers simultaneously.

```
memtier -s nsl200 -p 11120 -c <connections> -t 2  
  --random-data  
  --key-minimum=100  
  --key-maximum=10000
```

The effect will be to generate random data with keys between the specified ranges and send requests to the server in two simultaneous threads.

2.1.1 Default Throughput vs Latency

Firstly, we are interested in the relationship between throughput and latency shown in Figure 2.1. The mean latency and the 99th percentile latency are plotted against corresponding number of operations (throughput). We can observe that as the number of operations increases so does latency. Additionally, latency (both 99th percentile and mean) increase linearly until a saturation point is reached when a further increase in throughput is met with an exponentially larger increase in latency. The highest throughput achieved under quality of service restriction of 99th percentile latency under 1 millisecond is 375,000 operations per second. The highest level of throughput corresponds to 84 simultaneous connections, or 12 connections per each client which is similar to benchmarks used in the literature [8].

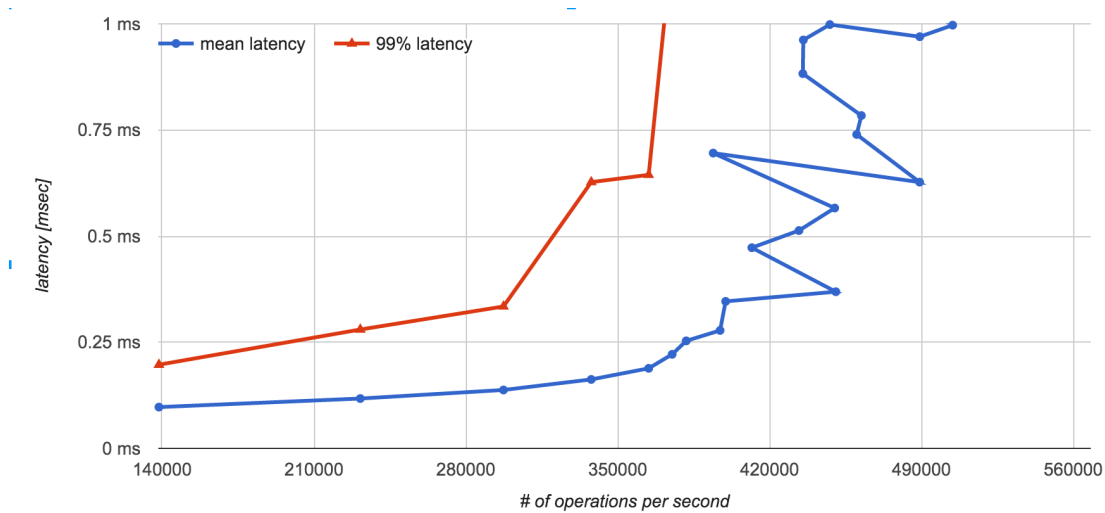


Figure 2.1: Mean latency and 99th percentile latency against throughput

2.1.2 Effect of Connections

To understand the effect of a large number of connections on the cache performance, Figure 2.2 shows the effect an increase in the total number of connections has on throughput, mean latency and the 99th percentile latency. The figure deliberately shows the behavior outside of the requires QoS requirements in order to better illustrate the impact on the cache under high load.

Firstly, Figure 2.2 displays the general trend an increased load has on throughput. As load increases, so does throughput. However, as the number of connections surpasses 100, the rate of increase in throughput for each increase the number of connections decreases. This is the saturation point of the cache, an increase in load yield disproportionate increase in throughput. Beyond the saturation point, the cache throughput fluctuates around 450,000 operations per second.

Secondly, the mean latency increases linearly with the number of connections (load). This is an expected behavior as the server experiences network stack queuing as well

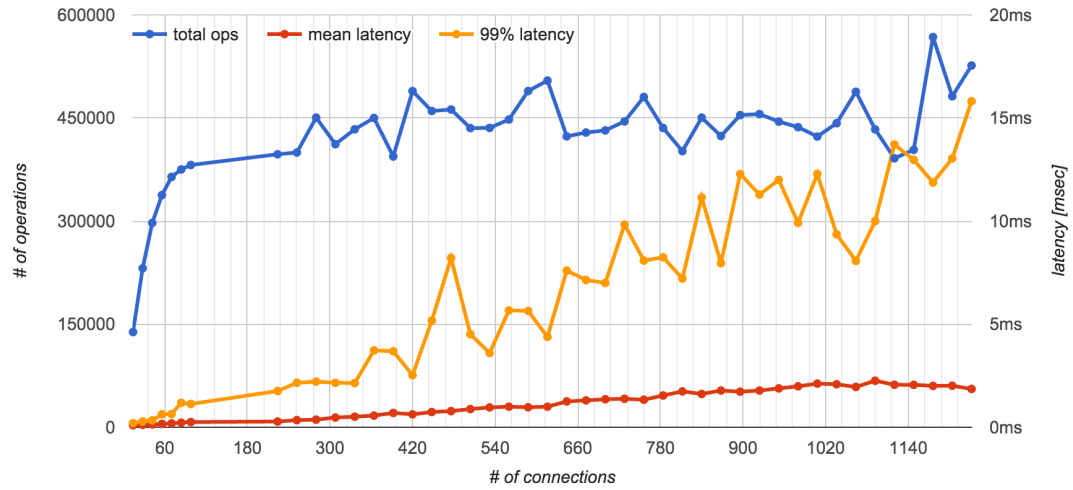


Figure 2.2: Throughput, Mean latency and 99th percentile against the number of connections

as increased resource requirements to process requests.

Thirdly, the 99th percentile latency increases linearly, with some fluctuations, against the increased server load. As the load gets higher, the fluctuation increases as well as the upper bound. This is due to some requests being queued for a long time before being processed, pushing the 99th percentile high.

We can observe that the server is capable of scaling much better until around 100 connections are reached. When the saturation point is surpassed, overall performance and quality of service degrades.

2.1.3 Server CPU

In order to be better understand the impact of memcached on the system, specifically the CPU usage, we consider a breakdown of CPU time spent in various areas of the operating system in Figure 2.3.

Firstly, we can observe that the effective footprint of memcached (*guest*) is relatively small compared to the rest. The CPU usage of memcached increases up until 100 connections at which point it remains fairly stable.

Secondly, the operating system (*sys*) increases as we increase the load. Fluctuations occur past 250 connections but the CPU usage by the system remains around 40 per cent as the load is increased further. An increased number of connections requires additional resources to process incoming requests as well as process outgoing requests. The context switching from receiving and transmitting contributes to the high load from the system.

Thirdly, interrupt processing by the system (*irq*) decreases as we increase load, this is due to having more CPU available and therefore being able to process a higher number of interrupts per unit of time. As resources are required by the system and the

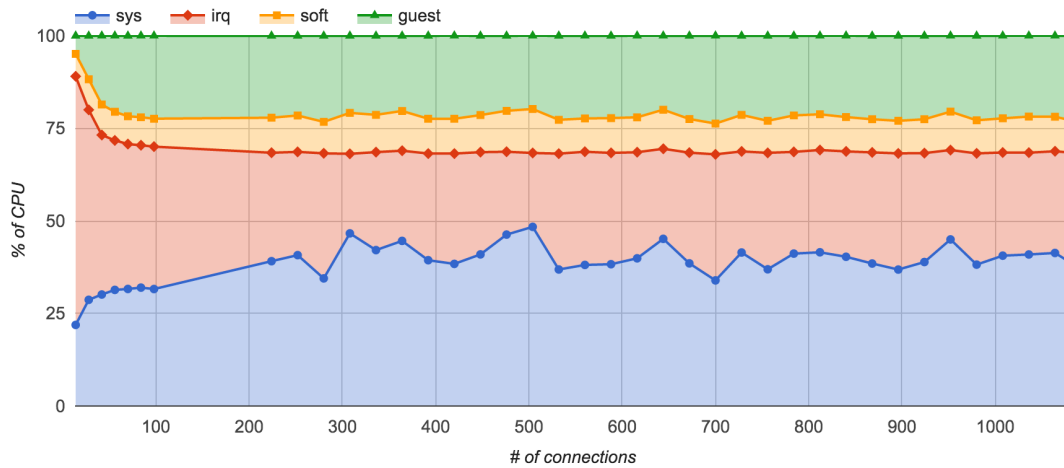


Figure 2.3: CPU Time against number of connections

memcached application, this number of interrupts processed per unit of time decreases as reflected by the proportion of CPU.

Finally, software interrupt footprint (*soft*) remains relatively stable throughout. The software interrupts correspond to running threads of the memcached application (4 threads by default) and are used for context switching.

From the breakdown in Figure 2.3, we can conclude that memcached is not CPU heavy on its own. The large CPU footprint of running memcached under high load is tightly linked to performance of the network stack and the underlying hardware processing network requests rather than the application itself. This observation is consistent with findings in MICA [7].

2.2 Memcached Thread Scalability

Memcached, as a high performance object cache, is designed to be executed on a parallel architecture. It implements scalability through the use multiple threads allowing memcached to utilize many core architectures. Therefore, the next step in scaling a memcached deployment is to provision a larger number of threads for the application.

Memcached execution model is capable of processing incoming and outgoing requests in parallel, however, operations executed require a global application lock to be acquired. Therefore, the expected number of threads maximising throughput while minimizing latency can be expected to be achieved when memcached is provisioned with the same number of threads as hardware CPU cores which is also suggested by Leverich and Kozyrakis [6].

Utilizing findings from the previous section, a configuration with 84 connections can be used to generate a consistent load while the number of threads provisioned for memcached can be varied. Therefore, we can set up each benchmark client as follows:

```
memtier -s nsl200 -p 11120 -c 6 -t 2 -P memcache_binary
```

```
--random-data
--key-minimum=100
--key-maximum=10000
```

The server in turn is configured as follows:

```
memcached -d -p 11120 -m 6144 -t <thread_count>
```

Where the number of threads is progressively increased.

2.2.1 Throughput & Latency

Figure 2.4 shows the relationship between throughput, mean latency and 99th percentile latency in relation to the number of threads used by a memcached application.

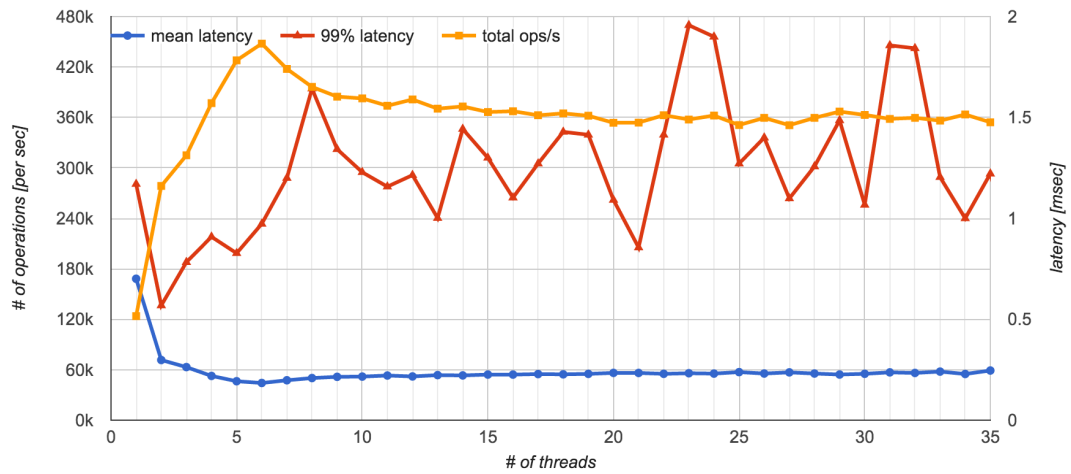


Figure 2.4: Memcached Thread Scaling

Firstly, we can observe that throughput increases as thread count increases until we reach 6 threads where it peaks at 450k requests per second. As we increase thread count further, throughput decreases. This behavior corresponds with our expectation that performance is maximized when there are as many threads as CPU cores.

Secondly, mean latency decreases as the number of threads is increased reaching a minimum of 0.1843 milliseconds at 6 threads. With more threads, the mean latency increases steadily.

Thirdly, the 99th percentile latency decreases as we increase the number of threads from 1 to 2, reaching a minimum and increasing as the number of threads increases. At 6 threads, we reach a 99th percentile latency of 0.973 which satisfies the QoS requirements under 1 millisecond.

Indeed, as expected we have been able to obtain the highest throughput and achieve the quality of service requirements with 6 threads, as many as CPU cores on the host. Beyond 6 threads, the overhead of context switching between threads increases processing time and reduces throughput.

2.2.2 CPU Time

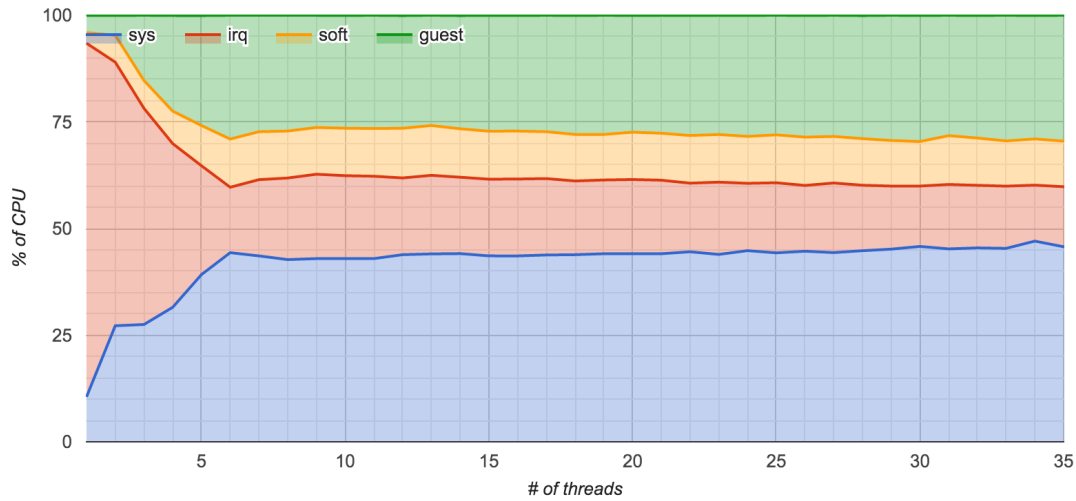


Figure 2.5: Memcached CPU Time against Number of Threads

Analyzing the CPU usage in Figure 2.5 as we increase the number of threads, we can observe that initially a large portion of the CPU time is spent servicing hardware interrupts (*irq*). Therefore, the OS is handling incoming traffic interrupts from the NIC. As the number of threads increases, an increasingly larger portion of CPU time is spent processing system calls and context switching (*sys*). This is reasonable as a larger number of threads will require context switching and concurrency management provided by the operating system. We can see that time spent processing hardware interrupts (*irq*) decreases which has the effect of increasing latency as packets remained queued up in the NIC for longer before the OS manages to schedule the interrupt to be serviced. Furthermore, we can observe that software interrupts (*soft*) CPU time progressively increases until we reach 6 threads and remains stable as the number of threads grows further. The initial increase is reasonable as we are demanding more threads to processed simultaneously, past this point the percentage remains stable as we have reached a saturation point in terms of scalability and server performance. Finally, memcached (*guest*) follows a similar pattern as software interrupts. Usage increases until 6 threads are used and saturates further. This is further indicative of the inability to efficiently scale the number of threads past the point at which memcached uses the same number of threads as CPU cores.

2.2.3 Thread evaluation

Comparing results obtain from thread scalability with the results from the default configuration of memcached, we have been able to increase throughput from 375k to 450k requests per second while maintaining the desired QoS under 1ms.

2.3 Thread pinning

Thread pinning is the process of assigning a *set_irq_affinity* to each individual thread. As suggested by Leverich and Kozyrakis, "pinning memcached threads to distinct cores greatly improves load balance, consequently improving tail latency." [6] and therefore the reasonable next step in optimizing memcached performance is to attempt thread pinning and analyse the results obtained.

By default, when a new process is started, its affinity is set to all available CPUs. We can discover a given process affinity by executing the following command where *pid* is the process identifier.

```
taskset -p <pid>
```

"A Memcache instance started with *n* threads will spawn *n* + 1 threads of which the first *n* are worker threads and the last is a maintenance thread used for hash table expansion under high load factor." [3]. We can discover memcached threads used for request processing using the following command where *tid* is the thread id discovered previously [3].

```
ps -p <memcache-pid> -o tid= -L | sort -n | tail -n +2 | head -n -1
```

Given the best performance under QoS constraints of 1ms found in the previous section is memcached with 6 threads, the following benchmark will be using this best configuration in order to analyze the impact of thread pinning.

2.3.1 Throughput vs Latency

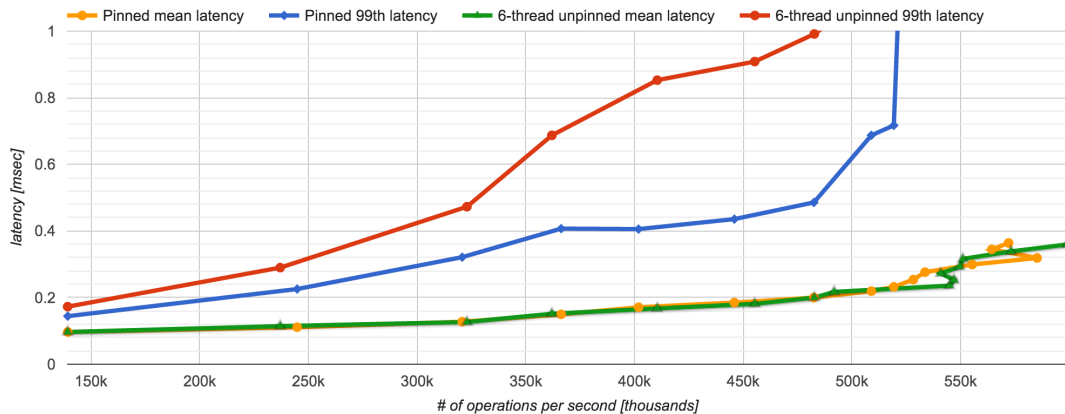


Figure 2.6: Memcached Pinned Threads vs Unpinned

Figure 2.6 shows the impact thread pinning has on mean and 99th percentile latency against throughput.

Firstly, the mean latency of both pinned and unpinned benchmarks remains very similar as throughput increases.

Secondly, the 99th percentile latency is lower in the case of pinned threads than unpinned. The pattern holds as throughput increases up until the required QoS boundary.

Furthermore, the throughput has also increased reaching 520k requests per second compared to 475k requests per second in the case of unpinned threads. This pattern is further confirmed by Leverich and Kozyrakis [6].

The reason for a significant improvement in the tail latency is reduced contention for access to receive and transmission queues in the network stack. Given memcached's design of running t threads for request processing and 1 thread for hash table expansion under high load, the result may be two or more request processing threads scheduled on the same CPU core while the hash table expansion thread is scheduled to run alone on a core. Pinning threads prevents this scenario from occurring as well as creates a direct mapping between each receive and transmission queue in the network stack to the memcached thread, reducing need for context switching.

2.4 Group Size

Memcached provides a configuration option `-R` to set the group size used inside memcached. The group size defines the “Maximum number of requests per event, limits the number of requests process for a given connection to prevent starvation (default: 20)” [4]. This in effect means the number of requests that will be processed from a single connection before memcached switches to a different connection to achieve a fair policy.

In order to benchmark the effects of an increased group size, we consider a benchmark scenario where the group size is progressively increased while maintaining a consistent load on the cache. The load used will be the one explored in Section 2.3 - Thread Pinning. All clients will be setup with 2 threads and 9 connections each which corresponds to the best performance under QoS so far. The memcached server will be setup with the `-R` flag set to 20 initially and increased by 20 for each iteration until a maximum (enforced by memcached) of 320 is reached.

Figure 2.7 shows the relationship between throughput, mean latency and 99th percentile latency.

Firstly, we can observe that the throughput achieved fluctuates heavily around 540k requests per second. Secondly, the mean latency remains constant as the group size is increased. Finally, the 99th percentile latency experiences a decreasing trend as group size is increased while staying under the required QoS constraints of 1ms.

By setting *group size* to 320, we have been able to reduce both 99th percentile latency as well as increase the total throughput from 520k to 540k requests per second. Additionally, we have been able to show that the context switching required to achieve a fairness policy of processing only 20 requests from a single connection caused a decrease in performance. Similar results have been obtained by Blake and Saidi [2]. On the other hand, the results observed are dependent on the distributed system model and

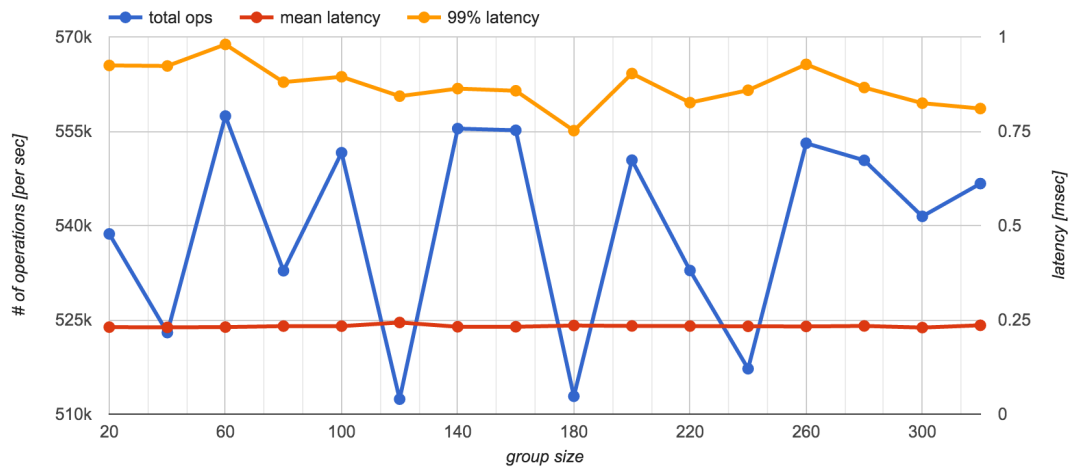


Figure 2.7: Memcached Latency, 99th percentile latency and throughput against group size

the expected request rate. In a production environment, benchmarking of a particular workload would show the best configuration in respect to group size.

2.5 Receive & Transmit Queues fixing

TODO: Haven't been able to obtain any improvements in terms of performance (both throughput and latency), not sure if the topic should still be discussed in detail.

2.6 Multiple Memcached Processes

In this section, we will examine the impact multiple memcached processes have on the overall performance of the caches as a whole. In some applications, it is important to be able to partition the system in such a way systems interact with different instances of memcached. Additionally, this information will also serve as a useful benchmark comparison for Redis performance.

Chapter 3

Redis

In this chapter, Redis performance in terms of latency, throughput and system resource requirements is examined. Initially, we will focus on performance under the default configuration of Redis. Subsequently, the scalability characteristics of Redis are explored. Focus is given to the impact of multiple Redis instances running simultaneously on the same machine under various levels of workload.

Unless otherwise stated, all benchmarks are performed to target the required Quality of Service (QoS) of achieving 99th percentile latency under 1 millisecond.

3.1 Out of the Box Performance

Firstly, in order to understand the baseline performance of Redis we consider the default configuration of Redis. A Redis deployment can be started with the following command:

```
redis redis.conf --port 11120
```

By default, a Redis deployment comes with a default configuration file `redis.conf`[10]. Any options specified in the configuration file can be overridden from the command line by prefixing them with `--`, in our case we are overriding the port number and setting it to 11120. All other configuration options remain unmodified.

In order to understand the default Redis performance, we design the benchmark to exert an increasing level of load on the Redis server. Initially, we start with 2 threads and 1 connection per each thread on all workload generating clients, there are 7 such clients in our setup. In our workload, we define the key space to be up to 100 million keys with an object size of 64 bytes yielding around 6.4 GB of data while the memory capacity of the server is 8GB and therefore data should not be swapped out of main memory. The workload generating clients are executed with the following command:

```
memtier -s nsl200 -p 11120  
        -c <connection_count> -t 2  
        -P redis
```

```
--random-data --data-size=64
--key-minimum=1 --key-maximum=100000000
--test-time=400
```

The number of connections per each thread is increased linearly in each subsequent benchmark run. Each workload is executed over the course of 400 seconds and the data is randomly generated with keys drawn from a uniform distribution between 1 and 100 million.

3.1.1 Latency vs Throughput

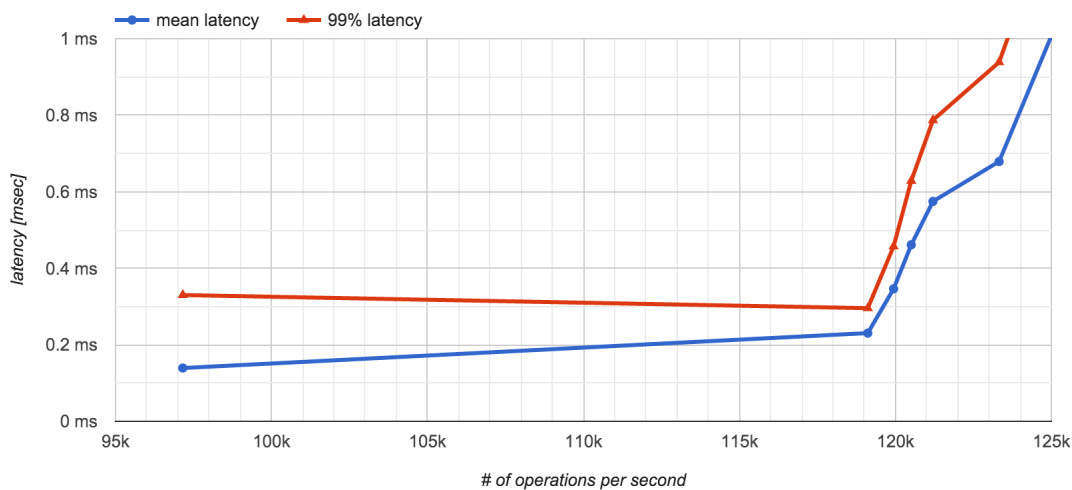


Figure 3.1: Redis: Throughput vs Mean and 99th percentile latency

Figure 3.1 plots the relationship between mean latency and the 99th percentile latency on the vertical axis and the number of operations per second on the horizontal axis. The graph has been trimmed to show only data which satisfies the QoS requirements.

We can observe that the number of operations Redis processes increases steadily until it reaches 119k requests per second at which point a further increase in throughput comes at a disproportionately greater cost in both mean and 99th percentile latency. The peak throughput observed under the QoS requirements is around 123k requests per second.

The performance degrades significantly with an increased workload beyond the peak shown at 123k requests per second. The throughput remains the same while the mean and the 99th percentile latency spikes heavily.

3.1.2 CPU Utilization

Figure 3.2 outlines the CPU utilization in terms of time spent processing system calls (sys), servicing hardware interrupts (irq), handling software interrupts (soft) and processing related to Redis (guest).

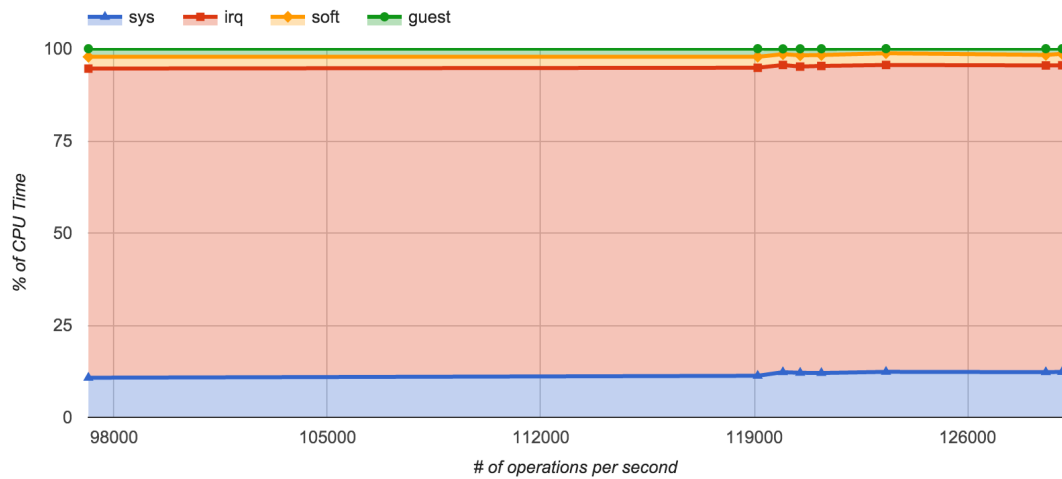


Figure 3.2: Redis: CPU Utilization

We can observe that servicing hardware interrupts consumes the majority of the processing time - 83 percent. This is the result of receiving and dispatching a large number of requests through the NIC. The amount of CPU time devoted to processing system calls is about 10 percent while processing software interrupts accounts for 3 percent. Redis itself requires only about 2 percent of the overall CPU time.

Redis appears to be network constrained rather than CPU constrained in the workload examined above. This is reasonable as Redis is executing within one main event loop therefore it does not require locking and does not generate heavy load on the CPU when looking values up in the cache.

3.2 Multiple Redis Instances

As seen in section 3.1, Redis cannot increase overall throughput of the system with only a single instance as only a single core is capable of processing requests and becomes the bottleneck. An immediate solution to this problem is to provision a larger number of instances on a multi-core system in order to better utilize the hardware resources.

In this benchmark, we examine the effects increased number of Redis instances with respect to latency, 99th percentile latency and overall throughput of the system. Additionally, we examine the effect of multiple instances on the CPU usage.

Firstly, multiple instances of Redis can be spawned easily on the server by binding them to distinct port numbers. Out of the box, Redis does not provide the capability to proxy multiple instances of Redis through a single port in order to load balance the instances. There is the option to configure a Redis cluster, however, the intended use case is primarily for resiliency and fail over. In this benchmark, we consider a simpler scenario where each instance is isolated from each other and acts as an independent cache. This is a simplification of a real world scenario, however, a large deployment

of Redis could be designed to partition the key space and utilize multiple independent instances similarly. The Redis application can be spawned with the following script:

```
for i in [1..n]
    redis-server redis.conf --port (11120 + i) --maxmemory (6 / i)gb
```

Note that we are explicitly specifying the maximum amount of memory each instance will be allocated. In our case, we partition 6 GB of memory space evenly between the individual instances.

Secondly, in order to obtain comparable results, the load exerted on the Redis cache must remain constant. The load itself, however, needs to be partitioned across all of the instances of Redis evenly. In order to achieve this, each workload generating client spawns *i* instances of the benchmark and targets its respective Redis instance. We use the following script to start the workload generating clients:

```
for i in [1..n]
    memtier -s nsl200 -p <port>
            -c round(16 / i)
            -t 1
            -P redis
            --random-data --data-size=64
            --key-minimum=1 --key-maximum=round(100000000 / i)
            --test-time=400
```

Initially, we start with 16 connections and 1 thread. As we increase the number of instances, the number of connections goes down, however, a larger number of instances are deployed. Note that we are using a `round` function to ensure that the number of connections as well as the maximum key are integers. In order to smooth out load variance caused by integer divisibility, we consider two cases. One in which the `round` function is defined as the `ceiling` function and the other when it is defined as the `floor` function. The results of both types of the `round` function are then averaged. If a higher number of workload generating clients were available for the experimentation, the `round` approach would not be required.

Furthermore, the generated dataset is 6.4GB (100 million keys * 64 bytes of data). This is by design and leads to evictions in the cache as the size approaches the maximum.

3.2.1 Latency and Throughput

Figure 3.3 plots the relationship between the number of Redis instances running simultaneously on the cache server against the mean and 99th percentile latency on the left vertical axis and the total number of operations per second on the right axis. The black line positioned at 1 ms outlines the QoS target of the benchmark.

The mean latency (blue) decreases steadily as the number of Redis instances increases. At five Redis instances, the mean latency flattens out and remains stable as the number of instances is increased further.

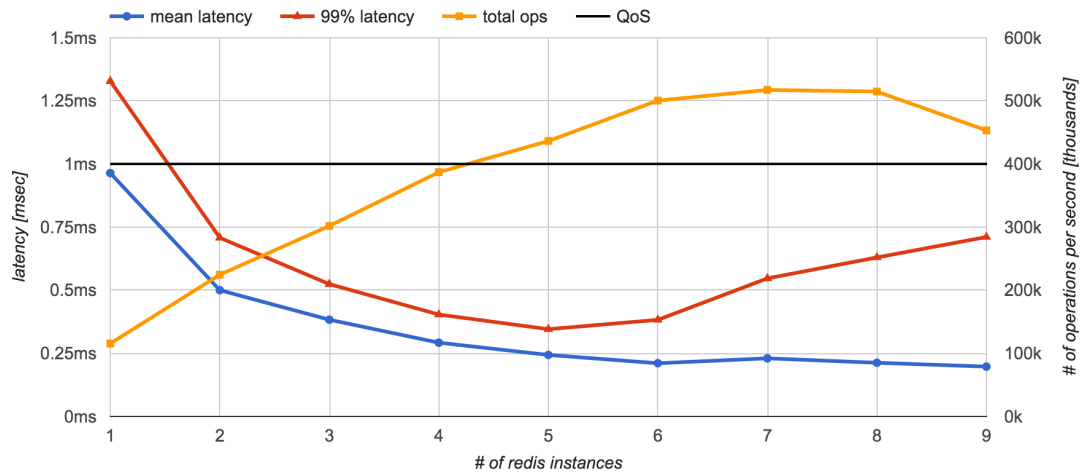


Figure 3.3: Redis Instances: Latency, Throughput vs Number of instances

The 99th percentile latency (red), on the other hand, experiences a sharp decrease as the number of instances is increased from 1 to 2 which brings the 99th percentile latency within the desired QoS constraints. As the number of instances increases further, the 99th percentile latency continues to decrease steadily. The minimum is reached at 5 instances with latency of 0.346 ms. A further increase in the number of instances results in higher 99th percentile latency.

The total number of operations increases steadily as the number of instances increases, peaking at 517,000 requests per second with 7 and 8 instances. A further increase in the number of instances results in a decrease in the total number of operations.

The QoS target is achieved in all but the 1 instance scenario. This is due to significantly larger load exerted on a single instance which in the subsequent multi-instance benchmarks gets distributed across the larger number of instances and therefore satisfies the QoS constraints.

Additionally, Redis manages to scale quite well as the number of instances increases up to 6. At this point, there are as many instances as there are CPU cores on the Redis server. Consequently, at 6 instances we achieve near minimum 99th percentile latency while also achieving near maximum throughput.

Furthermore, at 6 instances the 99th percentile latency is at around one third of the accepted QoS constraint. In real world deployment, it should be possible to utilize this gap to increase the throughput further. In our benchmarks, we have not been able to exert a load which would both increase throughput as well as increase the 99th percentile underutilization. This is due to a limited number of client benchmarking hosts to effectively tune the load to the required level.

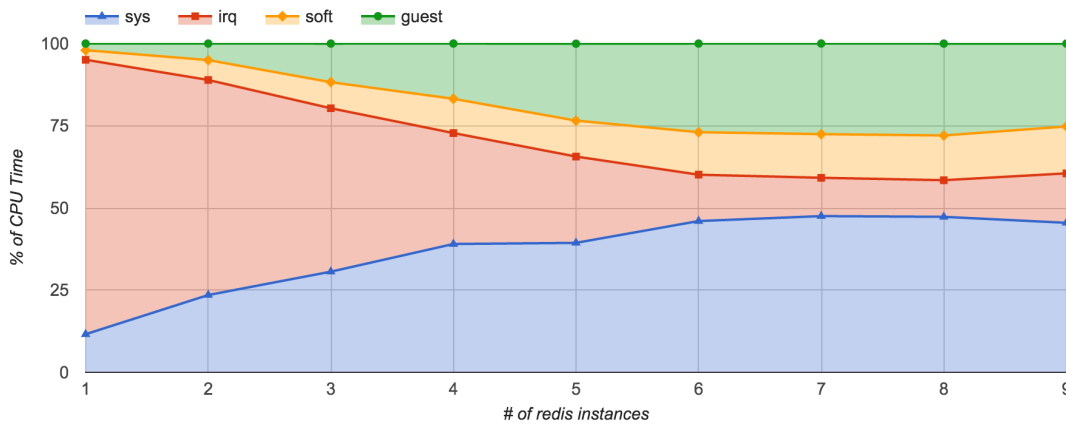


Figure 3.4: Redis Instances: CPU Utilization vs Number of Instances

3.2.2 CPU Utilization

Figure 3.4 presents the relationship between the number of instances and the CPU utilization.

Firstly, the CPU usage of Redis (green), *guest*, increases as the number of instances provisioned increases steadily. At 6 instances and more, the overall CPU usage remains fairly constant accounting for 27 percent of the total CPU utilization. This is due to the server having 6 CPU cores, requiring context switching between applications rather than parallel processing.

Secondly, the CPU dedicated to software interrupts (yellow), *soft*, follows a similar pattern to Redis. As the number of instances increases, so does the CPU usage dedicated to software interrupts. At 6 instances or more, the CPU utilization remains stable, accounting for 10 percent of the total usage. Software interrupts are inherently correlated to the software applications running on the server and therefore it is reasonable to observe a similar pattern between software interrupts and the application itself.

Thirdly, the CPU usage required by the operating system (blue), *sys*, grows steadily as the number of Redis instances increases. At 6 instances or more, the operating system CPU usage peaks and remains stable at 48 percent.

Finally, the portion of the CPU time dedicated to processing hardware interrupts (red), *irq*, decreases steadily as the number of Redis instances increases. This is due to batching of requests at the network interface card. At 6 or more instances, hardware interrupt servicing takes up about 11 percent of the total CPU usage.

Despite the increased resource requirements by Redis, it still only accounts for about 27 percent CPU usage at 6 processes with majority of the CPU time spent in the operating system and or networking. The observation from the single instance benchmark in the previous section therefore extends to a multi instance setup. Redis appears to be network intensive rather than CPU intensive.

By spawning multiple instances of Redis, we have been able to significantly increase

the overall performance of the server cache. Running multiple simultaneous instances, however, requires the key space to be partitioned.

3.3 Pinned Redis Instances

In the previous section we have observed that the performance of a Redis server can be greatly improved by provisioning multiple Redis instances simultaneously. Pinning processes to distinct cores is suggested to improve tail latency [6]. In this section, we examine the effect process pinning has on the performance of Redis. We consider exactly the same workload as in the previous section 3.2 as well as exactly the same server setup with the exception of pinning the Redis processes. That is, the workload is kept constant while it is partitioned across multiple instances.

A Redis process can be pinned to a unique core through the use of the `taskset` utility as follows:

```
taskset -pc <redis_pid> <core_id>
```

The Redis processes identified as `redis pid` is pinned to the CPU core identified by `core id`. We can identify the process id of a Redis application through the `ps` command. When running more Redis applications than there are CPUs, we assign it to the n th index of the application modulo the total number of cores, which is 6.

3.3.1 Pinned Latency and Throughput

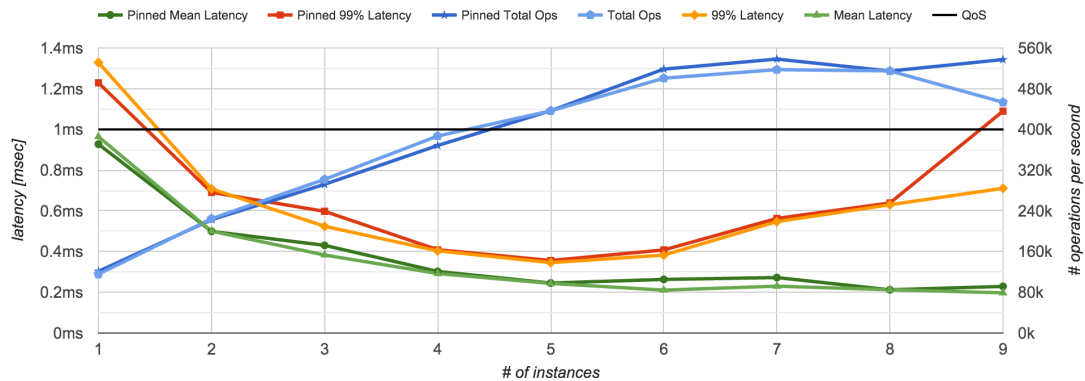


Figure 3.5: Redis Instance Pinning: Instances vs Latency and Throughput

Figure 3.5 plots the relationship between the number of instances on the horizontal axis, latency on the left vertical and throughput on the right vertical. Additionally, the performance obtained without pinning are plotted alongside the pinned results.

We can observe that pinning Redis processes results which strongly correlate to the results of the unpinned benchmark. Across mean latency, 99th percentile latency and throughput, there is very little variance in the performance observed.

3.3.2 Redis Persistence

TODO

3.4 Object Size

In this section, the impact of the size of the object stored in the cache is investigated. Redis imposes no restrictions on the size of objects stored in the cache.

In order to investigate the impact of object size on the cache, we consider a benchmark with an increasing object size. The object size is increased in powers of two starting at 2 bytes and ranging to 512 KB. This allows us to capture the majority of important sizes commonly used when designing applications.

The server configuration remains the same as the current best found configuration, the multi-instance configuration with 6 instances. The clients are configured as follows:

```
for i in [1..19]
  memtier -s nsl200 -p <port>
    -c 3
    -t 1
    -P redis
    --random-data --data-size=pow(2, i)
    --key-minimum=1 --key-maximum=(1066666666 / pow(2, i))
    --test-time=400
```

We run 19 iterations of the benchmark since 512 KB is equivalent to 2 to the power of 19 bytes. The data-size is configured to be increasing in powers of 2. The key range is defined as 6.4 GB split across 6 instances and further accounts for the increased size. Table 3.1 outlines the configuration options for each iteration.

3.4.1 Latency and Throughput

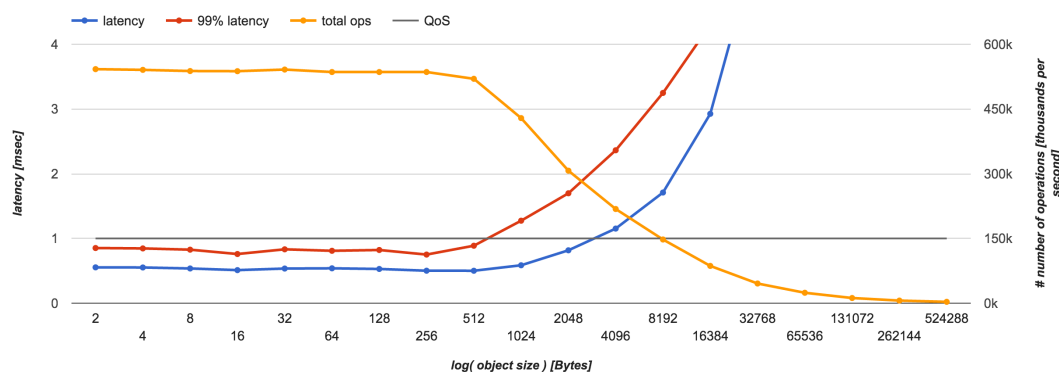


Figure 3.6: Redis Object Size: Latency and Throughput

Iteration	Data Size (bytes)	Key Maximum	Total Size (GB)
1	2	53333333	6.4
2	4	26666667	6.4
3	8	13333334	6.4
4	16	6666667	6.4
5	32	3333334	6.4
6	64	1666667	6.4
7	128	833334	6.4
8	256	416667	6.4
9	512	208334	6.4
10	1024	104167	6.4
11	2048	52084	6.4
12	4096	26042	6.4
13	8192	13021	6.4
14	16384	6511	6.4
15	32768	3256	6.4
16	65536	1628	6.4
17	131072	814	6.4
18	262144	407	6.4
19	524288	204	6.4

Table 3.1: Redis Object Size - Data Size and Maximum Key for each iteration. Total size is calculated as the product of Data Size, Key Maximum and 6 instances.

Figure 3.6 displays the relationship between object size on the horizontal logarithmic axis, latency on the left vertical axis and throughput on the right vertical axis.

Firstly, as object size increases up to 512 bytes, the mean latency remains stable at 0.55 ms. Beyond 512 bytes, the mean latency starts to increase and climbs beyond the desired QoS constraint at object size of 4 KB. A further increase in object size leads an disproportionately greater increase in mean latency.

Secondly, the 99th percentile follows the same pattern as the mean latency, however, it begins to climb over the desired QoS sooner, at object size of 1 KB.

Thirdly, the number of operations per second remains constant for object sizes under 256 bytes. An increase in object size decreases the number of operations per second. This is a reasonable result as an increase in the object size leads to higher bandwidth requirements and therefore leads to a lower number of operations per second.

Overall, Redis appears to be capable to scale well with object sizes up to 512 bytes. Larger object sizes put additional strain on the cache and require buffering which leads to increased latency of the average, and therefore 99th percentile, request. Primarily, Redis is not designed to store large (1KB+) values. It is, however, possible to partition large values into smaller ones and perform assembly/disassembly of the value on the client side.

3.4.2 CPU Utilization

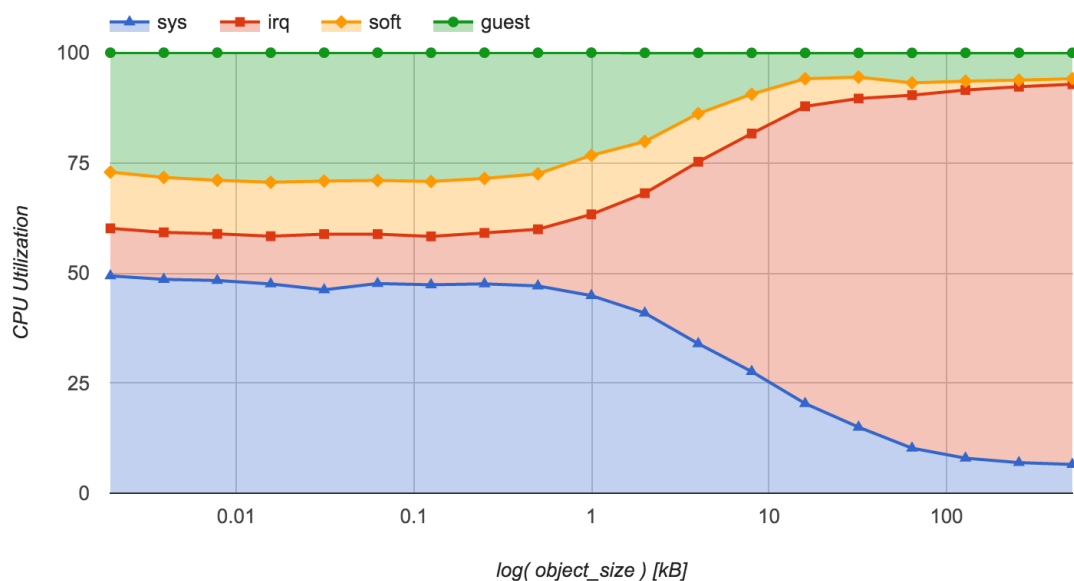


Figure 3.7: Redis Object Size: CPU Usage

Figure 3.7 displays the relationship between object size and CPU usage on the Redis server.

Firstly, as object size increases up to 1 KB, the operating system (sys) requires about 47 percent of the total time to process the incoming requests and dispatch them to the relevant applications. As object size increases further, the time required decreases as the number of operations decreases.

Secondly, the Redis applications (guest) require 27 percent of the total time when processing requests under 1 KB, with requests larger the direct cost of running Redis decreases as there are less requests to process. A similar pattern holds for the software interrupts (soft), as there are less requests coming.

Finally, the time allocated to servicing hardware interrupts (irq) remains at 12 percent below 1KB, with object size increases beyond 1 KB, there is significant increase in the time required to service hardware interrupts. This is due to buffering of large objects and is effectively the cause of high mean and 99th percentile latency as well as low throughput.

Overall, Redis is designed to work well with objects sizes below 1 KB. As the object size increases, the cache experiences a degraded performance due to buffering of network input and output.

3.5 Key Distributions

TODO

3.5.1 Gaussian distribution

TODO

3.5.2 Zipf distribution

TODO

Chapter 4

Redis & Memcached: Head to Tail

Evaluation goes here

Chapter 5

Conclusion

Bibliography

- [1] Berk Atikoglu, Yuehai Xu, Eitan Frachtenberg, Song Jiang, and Mike Paleczny. Workload analysis of a large-scale key-value store. In *ACM SIGMETRICS Performance Evaluation Review*, volume 40, pages 53–64. ACM, 2012.
- [2] Geoffrey Blake and Ali G Saidi. Where does the time go? characterizing tail latency in memcached. *System*, 54:53.
- [3] Solarflare Communications Inc. Filling the pipe: A guide to optimising memcache performance on solarflare hardware. 2013.
- [4] Danga Interactive. Memcached. <http://memcached.org>.
- [5] Redis Labs. Memtier benchmark. https://github.com/RedisLabs/memtier_benchmark.
- [6] Jacob Leverich and Christos Kozyrakis. Reconciling high server utilization and sub-millisecond quality-of-service. In *Proceedings of the Ninth European Conference on Computer Systems*, page 4. ACM, 2014.
- [7] Hyeontaek Lim, Dongsu Han, David G Andersen, and Michael Kaminsky. Mica: A holistic approach to fast in-memory key-value storage. *management*, 15(32):36, 2014.
- [8] Kevin Lim, David Meisner, Ali G Saidi, Parthasarathy Ranganathan, and Thomas F Wenisch. Thin servers with smart pipes: designing soc accelerators for memcached. *ACM SIGARCH Computer Architecture News*, 41(3):36–47, 2013.
- [9] Niels Provos and Nick Mathewson. libevent an event notification library. <http://libevent.org>.
- [10] Redis. redis.conf. <https://github.com/antirez/redis/blob/3.0/redis.conf>, 2015.
- [11] RedisLabs. memtier_benchmark: A high-throughput benchmarking tool for redis and memcached. https://redislabs.com/blog/memtier_benchmark-a-high-throughput-benchmarking-tool-for-redis-memcached#.VunpLhKLTmu, 2013.
- [12] Alex Wiggins and Jimmy Langston. Enhancing the scalability of memcached. *Intel document, unpublished*, 2012.