

RDMA for Memcached 0.9.6 User Guide

HIGH-PERFORMANCE BIG DATA TEAM
<http://hibd.cse.ohio-state.edu>

NETWORK-BASED COMPUTING LABORATORY
DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
THE OHIO STATE UNIVERSITY

Copyright (c) 2011-2017
Network-Based Computing Laboratory,
headed by Dr. D. K. Panda.
All rights reserved.

Last revised: December 23, 2017

Contents

| | | |
|----------|-----------------------------------------------------------------------------------|-----------|
| 1 | Overview of the RDMA for Memcached Project | 1 |
| 2 | Features | 1 |
| 3 | Installation Instructions | 2 |
| 3.1 | Pre-requisites | 2 |
| 3.2 | Download | 2 |
| 3.3 | Installing from RHEL6 package | 2 |
| 3.4 | Installing from tarball | 3 |
| 4 | Runtime Parameters | 3 |
| 4.1 | Memcached Server Parameters | 3 |
| 4.2 | Memcached Client Parameters | 4 |
| 4.2.1 | Enabling RoCE mode | 4 |
| 4.2.2 | Runtime IB/RoCE HCA selection | 4 |
| 5 | RDMA-Memcached Server | 4 |
| 5.1 | Server Runtime Modes | 5 |
| 5.2 | Running RDMA-Memcached Server with In-Memory Mode | 5 |
| 5.3 | Running RDMA-Memcached Server with Hybrid-Memory Mode | 5 |
| 5.4 | Running RDMA-Memcached Server with Burst-buffer Mode | 6 |
| 5.5 | Running RDMA-Memcached Server over RoCE clusters | 6 |
| 5.6 | Enable IB/RoCE HCA Selection at Runtime | 6 |
| 6 | RDMA-Memcached Client | 7 |
| 6.1 | Libmemcached Blocking API | 7 |
| 6.1.1 | Description | 7 |
| 6.1.2 | Client Example | 7 |
| 6.2 | Libmemcached Non-Blocking API for RDMA-Memcached | 8 |
| 6.2.1 | Description | 8 |
| 6.2.2 | Client Example | 10 |
| 6.3 | Libmemcached Support for RDMA-based HDFS (HHH) Burst-Buffer Mode | 10 |
| 6.4 | Enabling RDMA-Memcached Client to run over RoCE | 10 |
| 6.5 | Enable IB/RoCE HCA Selection at Runtime | 10 |
| 7 | Memcached Micro-benchmarks | 11 |
| 7.1 | OHB Memcached Micro-benchmarks | 11 |
| 7.1.1 | Memcached Latency Micro-benchmark (ohb_memlat) | 11 |
| 7.1.2 | Memcached Hybrid Micro-benchmark (ohb_memhybrid) | 12 |
| 7.1.3 | Memcached Latency Micro-benchmark for Non-Blocking APIs (ohb_memlat_nb) | 13 |
| 7.2 | YCSB extension for RDMA-Memcached in OHB | 15 |
| 7.2.1 | Building YCSB with RDMA-Libmemcached Support | 15 |
| 7.2.2 | Running the YCSB Benchmark Extension with RDMA-Memcached | 16 |
| 8 | Troubleshooting with RDMA-Memcached | 17 |

1 Overview of the RDMA for Memcached Project

RDMA for Memcached is a high-performance design of Memcached over RDMA-enabled Interconnects. In addition to enabling low latencies on InfiniBand and 10/40GigE cluster, we introduce new features into the Memcached design including support for high-performance HDFS (HHH) burst-buffer mode and non-blocking Libmemcached API semantics. This version of RDMA for Memcached 0.9.6 is based on Memcached 1.5.3, and Libmemcached 1.0.18. This file is intended to guide users through the various steps involved in installing, configuring, and running RDMA for Memcached over InfiniBand. This guide provides examples to illustrate the use of both traditional blocking Set/Get APIs (`memcached_set/memcached_get`) and newly introduced non-blocking Set/Get APIs (`memcached_liset/memcached_liget/memcached_bset/memcached_bget`). This guide also describes the micro-benchmarks defined for experimentation with the different advanced features in the RDMA for Memcached package.

If there are any questions, comments or feedbacks regarding this software package, please post them to `rdma-memcached-discuss` mailing list (`rdma-memcached-discuss@cse.ohio-state.edu`).

2 Features

High-level features of RDMA for Memcached 0.9.6 are listed below. New features and enhancements compared to 0.9.5 release are marked as **(NEW)**.

- **(NEW)** Based on Memcached 1.5.3
 - **(NEW)** Compliant with the Memcached's new item chaining feature in In-Memory mode
 - **(NEW)** Compliant with the latest Memcached's LRU maintainer and slab balancer enhancements
- Based on libMemcached 1.0.18
- High performance design with native InfiniBand and RoCE support at the verbs level for Memcached Server and Client
- High performance design of SSD-assisted hybrid memory
- **(NEW)** Runtime selection of HCA device for nodes equipped with multiple InfiniBand/RocE HCAs
- **(NEW)** Enable and disable item chaining through extended server options
- Compliant with libMemcached 1.0.18 APIs and applications
- Non-Blocking Libmemcached Set/Get API extensions
 - APIs to issue non-blocking set/get requests to the RDMA-based Memcached servers
 - APIs to support monitoring the progress of non-blocking requests issued in an asynchronous fashion
 - Facilitating overlap of concurrent set/get requests
- Support for burst-buffer mode in Lustre-integrated design of HDFS in RDMA for Apache Hadoop-2.x

- Support for both RDMA-enhanced and socket-based Memcached clients
- Easily configurable for native InfiniBand, RoCE, and the traditional sockets based support (Ethernet and InfiniBand with IPoIB)
- On-demand connection setup
- Tested with
 - Native Verbs-level support with Mellanox InfiniBand adapters (QDR, FDR, and EDR)
 - RoCE support with Mellanox adapters
 - Various multi-core platforms
 - SATA-SSD, PCIe-SSD, and NVMe-SSD

3 Installation Instructions

3.1 Pre-requisites

In order to use the RDMA-based features provided with RDMA for Memcached, install the latest version of the OFED distribution that can be obtained from <http://www.openfabrics.org>. It also requires cyrus-sasl-devel package (eg: `yum install cyrus-sasl-devel`) and the libhugetlbfs package (<http://libhugetlbfs.sourceforge.net>).

3.2 Download

The RDMA for Memcached package consists of Memcached Server executable and Libmemcached Client libraries.

The rpm for the latest version of RDMA-Memcached package can be downloaded from:

```
http://hibd.cse.ohio-state.edu/download/hibd/rdma-memcached-0.9.6-1.el6.x86_64.rpm.
```

The latest version of RDMA-Memcached is also available as a tarball at:

```
http://hibd.cse.ohio-state.edu/download/hibd/rdma-memcached-0.9.6-x86-bin.tar.gz.
```

3.3 Installing from RHEL6 package

Running the following command script will install the software in `/usr/local/rdma-memcached-0.9.6`.

```
rpm -Uvh rdma-memcached-0.9.6-1.el6.x86_64.rpm
```

Alternatively, users can use `--prefix=PATH` to install in specific path.

The above command will upgrade any prior version of RDMA-Memcached that may be present.

3.4 Installing from tarball

To install using tarball,

- Extract memcached and libmemcached binaries from the tarball downloaded:

```
tar -xf rdma-memcached-0.9.6-x86-bin.tar.gz
```
- Change directory to find RDMA-based memcached binary and libmemcached libraries:

```
cd rdma-memcached-0.9.6
```

Please email us at rdma-memcached-discuss@cse.ohio-state.edu if your distro does not appear on the list or if you experience any trouble installing the package on your system.

4 Runtime Parameters

Some advanced features in RDMA for Memcached 0.9.6 can be manually enabled by users. Runtime parameters in RDMA for Memcached 0.9.6 include:

4.1 Memcached Server Parameters

We describe the runtime parameters that are specific to RDMA-Memcached Server in this section.

- p** **<num>**: TCP port to listen on for RDMA connections (default: 11211)
- t** **<num>**: Number of threads to use to handle RDMA clients (default: 4)
- m** **<num>**: Maximum item memory in megabytes (default: 64 MB). For hybrid mode, this represents key/index memory.
- H** **<SSD-path>**: Enable SSD-assisted hybrid mode with RDMA Memcached. This is used to specify a fully qualified file-path. File-path must specify an existing directory on SSD.
- z**: Enable direct I/O for SSD read/write (O_DIRECT flag)
- N** **<num>**: IB/RoCE HCA device number to use (default: 0)
- W** **<num>**: IB/RoCE HCA device name to use
- e** **<num>**: Value item memory in megabytes (default: 64 MB) This valid only in hybrid mode.
- g** **<num>**: Maximum SSD file size in megabytes (default: 32 GB) This is applicable to the RDMA Memcached hybrid mode.
- T** **<num>**: number of threads to use to handle RDMA clients (default: 0)
- E** **<num>**: TCP port to listen on (default: 11212, not enabled by default)

- l** **<addr>**: Interface to listen on (default: INADDR_ANY, all addresses) **<addr>** may be specified as host:port. If you don't specify a port number, the value you specified with **-p** or **-U** is used. You may specify multiple addresses separated by comma or by using **-l** multiple times
- d**: Run as a daemon
- h**: Print this help and exit
- i**: Print memcached and libevent license
- o** **chunked_items**: Enable item chaining explicitly for RDMA-Memcached In-Memory mode. This option is not valid for hybrid mode or burst-buffer mode.
- o** **maxbuffers_chunked=<num>**: Maximum number of temporary buffers to use for remote fetches and per thread with item chaining i.e., **-o** **chunked_items**. (default: 8 and fixed-size of **val_size_max**).
- o** **use_roce**: Use RDMA-over-Converged-Ethernet (RoCE) mode for RDMA workers.
- o** **burst_buffer=<SSD-dir-path>**: Enable burst-buffer mode for RDMA-Hadoop HDFS (HHH-L-BB). Provide a valid SSD directory path for persisting data. i.e., **-o** **burst_buffer=<SSD-dir-path>**.

4.2 Memcached Client Parameters

We describe the runtime parameters that are specific to RDMA-Memcached Client in this section.

4.2.1 Enabling RoCE mode

To enable the RDMA-Memcached client to run over RDMA-over-Converged-Ethernet (RoCE), please set the environment variable **MEMCACHED_USE_ROCE** to 1. An example is provided in Section 6.4.

4.2.2 Runtime IB/RoCE HCA selection

The RDMA-Memcached client can choose the InfiniBand or RoCE HCA on a multi-HCA node by setting either of the two environment variables: **MEMCACHED_USE_HCA_NUM** or **MEMCACHED_USE_HCA_NAME**. Examples are provided in Section 6.5.

5 RDMA-Memcached Server

Memcached Server can be setup and run similar to the sockets-based Memcached server. For a list of arguments, please run **<RDMA-MEMCACHED_INSTALL_PATH>/memcached/bin/memcached -h**.

5.1 Server Runtime Modes

Our package can supports two primary modes: In-Memory and Hybrid-Memory modes. In addition to this, we support the Memcached-based burst-buffer mode for high-performance HDFS available in the RDMA for Hadoop 1.3.0.

- **In-memory mode:** This mode adheres to the default design of memory management in Memcached. When memory is limited, the memcached server either evicts older key/value pairs to make place for new key/value pairs, or returns Out-Of-Memory error.
In this latest release, item chaining introduced in default Memcached version 1.4.29 can be enable for the in-memory mode. With item chaining, memory efficiency can be enabled by re-using memory from smaller slabs for a given item. Details about item chaining can be found [here](#).
- **Hybrid-memory mode:** This mode extends the memory management schemes in Memcached to use SSD's assistance when memory is limited, rather that discarding or erroring out, in order to accommodate more key/value pairs in Memcached.
- **Burst-buffer mode for HHH:** This mode enables high-performance RDMA-based HDFS (HHH) to use Memcached server as a burst-buffer mode (HHH-L-BB). This mode uses Memcached server cluster to alleviate the load on the Lustre parallel file system while using the HHH-L mode and takes advantage of high-speed SSDs local to the Memcached servers for fault-tolerance.

5.2 Running RDMA-Memcached Server with In-Memory Mode

For default in-memory mode, please run

```
$$> <RDMA-MEMCACHED_INSTALL_PATH>/memcached/bin/memcached -t  
    <num-threads> -m <max-memory-in-megabytes>
```

To enabled RDMA-Memcached server to chunk-and-stitch items for item chaining, we buffer remote RDMA reads at the server threads. The number of such buffers can be specified, especially to help handle non-blocking set or get requests.

```
$$> <RDMA-MEMCACHED_INSTALL_PATH>/memcached/bin/memcached -t  
    <num-threads> -m <max-memory-in-megabytes>  
    -o chunked_items,maxbuffers_chunked=<max-buffers-per-thread>
```

5.3 Running RDMA-Memcached Server with Hybrid-Memory Mode

For hybrid-memory mode, please run

```
$$> <RDMA-MEMCACHED_INSTALL_PATH>/memcached/bin/memcached -t  
    <num-threads> -H <path-to-file-on-SSD>
```

```
-m <max-memory-for-keys-in-megabytes>
-e <max-memory-for-values-in-megabytes>
-g <max-SSD-file-size-limit-in-megabytes>
```

Normal Cached I/O is used to access SSD in the hybrid mode. If desired, direct I/O can be enabled using `-z` runtime parameter.

5.4 Running RDMA-Memcached Server with Burst-buffer Mode

For burst-buffer mode, please run

```
$$> <RDMA-MEMCACHED_INSTALL_PATH>/memcached/bin/memcached -t
<num-threads> -m <max-memory-for-keys-in-megabytes>
-o burst_buffer=<ssd-dir-path-to-persist-hdfs-data>
```

This runtime mode is introduced to support the burst-buffer mode of HHH design in RDMA for Hadoop 1.3.0 that leverages both local storage and Lustre for HDFS I/O. Please see [here](#) for more information on the RDMA for Hadoop package.

5.5 Running RDMA-Memcached Server over RoCE clusters

To enable the RDMA-Memcached server to run over RDMA-over-Converged-Ethernet (RoCE), the extended server option `use_roce` can be used. This can be enabled for all three modes mentioned above. For example, to enable RoCE with In-Memory mode:

```
$$> <RDMA-MEMCACHED_INSTALL_PATH>/memcached/bin/memcached -t
<num-threads> -m <max-memory-in-megabytes>
-o use_roce,<other-extended-options>
```

5.6 Enable IB/RoCE HCA Selection at Runtime

The RDMA-Memcached server can choose the InfiniBand/RoCE HCA for either of the three modes (In-Memory, Hybrid-Memory, Burst-Buffer) using the `-N` or `-W` server options, on nodes equipped with multiple InfiniBand and/or RocE HCAs.

For instance, to choose the HCA by device number for In-Memory mode:

```
$$> <RDMA-MEMCACHED_INSTALL_PATH>/memcached/bin/memcached -t
<num-threads> -m <max-memory-in-megabytes>
-N <hca-device-number>
```

Similarly, to choose HCA by name for the In-Memory Mode:


```
$$> <RDMA-MEMCACHED_INSTALL_PATH>/memcached/bin/memcached -t  
    <num-threads> -m <max-memory-in-megabytes>  
    -W <hca-device-name>
```

Information regarding the HCAs installed on a given node, i.e., device number, device name, and their corresponding transport protocol e.g, IB or Ethernet (RoCE over 40Gbe), the `ibstat` command can be run.

6 RDMA-Memcached Client

The RDMA for Memcached package supports APIs in the default `libmemcached` library for the Memcached binary protocol. In addition to these, we introduce RDMA-enabled non-blocking set and get APIs in RDMA for Memcached 0.9.6. These APIs enable the users to issue set or get requests without needing to block on the response. and the response of individual requests can be monitored in an asynchronous fashion.

6.1 Libmemcached Blocking API

6.1.1 Description

We support the default APIs with the blocking behaviour, including `memcached_set`, `memcached_get`, `memcached_mget`, `memcached_increment`, `memcached_decrement` and their variations. Their implementations are designed to be compatible with existing API semantics.

6.1.2 Client Example

Memcached client programs need to be compiled with RDMA for Memcached Client library. A simple Memcached client program is provided along with the package, and it gets installed in `<RDMA-MEMCACHED_INSTALL_PATH>/libmemcached/share/example` folder. This example program just issues a Set operation to the server, and then does a Get operation to retrieve the value from the server. The status of operations are output. The Memcached server node and port can be specified as command line argument.

Users can compile Memcached Client applications with RDMA for Memcached Client library.

```
$$> gcc memcached_example.c -o memcached_example  
    -I<RDMA-MEMCACHED_INSTALL_PATH>/libmemcached/include  
    -L<RDMA-MEMCACHED_INSTALL_PATH>/libmemcached/lib  
    -lmemcached -lpthread
```

The client application can be run be executed, and any environmental parameters can be specified at runtime.

```
$$> EXAMPLE_ENV_PARAM=VALUE ./memcached_example node123:11211  
$$> Key stored successfully  
$$> Key retrieved successfully. Key = KEY, Value = VALUE123
```

In this example, the client program does a Memcached Set operation to the server running on node123 at port 11211, followed by a Get operation from the same server.

6.2 Libmemcached Non-Blocking API for RDMA-Memcached

We introduce novel and high-performance non-blocking key/value store API semantics in RDMA for Memcached 0.9.6, that can co-exist with default blocking APIs. These non-blocking API extensions allow the user to separate the request issue and completion phases. The proposed APIs can issue set/get requests and return to the user application as soon as the underlying RDMA communication engine has communicated the request to the Memcached servers; without blocking on the response. Supporting progress APIs are available to check the progress of each operation in either a blocking or non-blocking fashion, respectively. We benefit from the inherent one-sided characteristics of the underlying RDMA communication engine, while ensuring the completion of every read or write request. By exploiting these non-blocking semantics to issue I/O requests, the client-side waits and server-side memory management and/or I/O can be overlapped with other data accesses or computations. For more information, please find our publication [here](#).

6.2.1 Description

Non-blocking Request Structures: To facilitate the proposed non-blocking API semantics, we introduce a new structure called `memcached_request_st` that contains: (1) completion flag that the user can test or wait on for operation completion, (2) pointer to the buffer where the server response will be available, and (3) pointers to user's value buffer. For each non-blocking request, a unique `memcached_request_st` is created and is used to monitor the receipt of response from the server in an asynchronous fashion. APIs for creating and destroying these request items are as follows:

```
memcached_request_st *memcached_request_create(memcached_st *ptr,
                                              void *value, size_t value_length, uint32_t flags);
```

This function is used to create a non-blocking request structure that will then be used by other non-blocking RDMA-enabled Libmemcached functions to communicate with the server. It takes `value` buffer pointer, `value` length and any flags to be used for the request. It returns a pointer to the `memcached_request_st` created or NULL on errors.

```
void memcached_request_free(memcached_st *ptr, memcached_request_st *req);
```

This function is used to clean up memory associated with a `memcached_request_st` structure. You should use this when you have received a response.

Non-blocking Set APIs: We introduce non-blocking APIs for set requests with and without buffer re-use guarantees. With buffer re-use guarantee, the `value` buffer used to create the `memcached_request_st` can be re-used or returned once the request has successfully completed. On the other hand, without buffer re-use guarantee, the `value` buffer cannot be re-used until a response has been received (via progress APIs).

```
memcached_return_t memcached_iset(memcached_st *ptr, const char *key,
```

```
size_t key_length, memcached_request_st *req,  
time_t expiration, uint32_t flags);
```

This function issues a set request without buffer re-use guarantees. It takes a key and its length, along with the `memcached_request_st` created. On success the value will be `MEMCACHED_SUCCESS`.

```
memcached_return_t memcached_bset(memcached_st *ptr, const char *key,  
size_t key_length, memcached_request_st *req,  
time_t expiration, uint32_t flags);
```

This function issues a set request with buffer re-use guarantees. It takes a key and its length, along with the `memcached_request_st` created. On success the value will be `MEMCACHED_SUCCESS`.

Use `memcached_strerror()` to translate returned status to a printable string.

Non-blocking Get APIs: Complimentary to set, we introduce non-blocking APIs for get requests with and without buffer re-use guarantees.

```
memcached_return_t memcached_iget(memcached_st *ptr, const char *key,  
size_t key_length, memcached_request_st *req);
```

This function issues a get request without buffer re-use guarantees. It takes a key and its length, along with the `memcached_request_st` created to receive the response and the value from the Memcached server. On success the value will be `MEMCACHED_SUCCESS`.

```
memcached_return_t memcached_bget(memcached_st *ptr, const char *key,  
size_t key_length, memcached_request_st *req);
```

This function issues a get request with buffer re-use guarantees. It takes a key and its length, along with the `memcached_request_st` created to receive the response and the value from the Memcached server. On success the value will be `MEMCACHED_SUCCESS`.

Use `memcached_strerror()` to translate returned status to a printable string.

Non-blocking Progress APIs: Non-blocking progress APIs check the progress of the operation in either a blocking or non-blocking fashion. The APIs are as follows:

```
void memcached_test(memcached_st *memc, memcached_request_st *req);
```

This function is used to monitor progress of a non-blocking set or get request in a non-blocking fashion. If response has been received it returns `MEMCACHED_SUCCESS` on success or `MEMCACHED_FAILURE`. If the response has not yet been received it returns `MEMCACHED_IN_PROGRESS`.

```
void memcached_wait(memcached_st *memc, memcached_request_st *req);
```

This function is used to monitor progress of a non-blocking set or get request in a blocking fashion. It returns `MEMCACHED_SUCCESS` on success request completion, else, it returns `MEMCACHED_FAILURE`.

6.2.2 Client Example

Memcached client programs using non-blocking APIs need to be compiled with RDMA for Memcached Client library. An example that illustrates how the `iset/bset/iget/bget` operations is used is provided along with the package, and it gets installed in `<RDMA-MEMCACHED-INSTALL_PATH>/libmemcached/share/example` folder. This example program (`memcached_nb_example.c`) can be compiled similar to the example described in Section 6.1.2.

6.3 Libmemcached Support for RDMA-based HDFS (HHH) Burst-Buffer Mode

We introduce configurable parameters to enable the burst-buffer mode for Hadoop users running applications with HHH and underlying file system with Lustre integrated support. This feature can be enabled in `hdfs-site.xml` by setting `hadoop.bb.enabled` to `true` and `memcached.server.list` with a comma-separated list of Memcached server hostnames.

6.4 Enabling RDMA-Memcached Client to run over RoCE

In this section, we provide an example for running RDMA-Memcached Client running in RoCE mode. Ensure that `-o use_roce` is specified at the RDMA-Memcached servers.

For instance, to run `memcached_example` with RoCE enabled:

```
$$> MEMCACHED_USE_ROCE=1 ./memcached_example node123:11211
$$> Key stored successfully
$$> Key retrieved successfully. Key = KEY, Value = VALUE123
```

6.5 Enable IB/RoCE HCA Selection at Runtime

In this section, we provide an example for choosing a particular IB/RoCE HCA by HCA device name or device number. Ensure that appropriate HCA is chosen at the RDMA-Memcached servers using `-N` or `-W` option.

For instance, to choose the HCA with device number 1 and name 'mlx4_0', either of the two ways can be used to run the client:

```
$$> MEMCACHED_USE_HCA_NUM=1 ./memcached_example node123:11211
```

OR

```
$$> MEMCACHED_USE_HCA_NAME="mlx4_0" ./memcached_example node123:11211
```

7 Memcached Micro-benchmarks

The OHB Micro-benchmarks support stand-alone evaluations of Memcached, Hadoop Distributed File System (HDFS), HBase and Spark (See [here](#)). OSU HiBD-Benchmarks (version 0.9.3) for Memcached consist of Get and Set latency micro-benchmarks.

7.1 OHB Memcached Micro-benchmarks

The source code can be downloaded from [osu-hibd-benchmarks-0.9.3.tar.gz](#). The source can be compiled with the help of the Maven (version 3.3.0 or higher) as follows:

1. Ensure that `libmemcached.home` property is set to the fully-qualified RDMA for Memcached 0.9.6 install path in `OHB_INSTALL_PATH/memcached/microbench/pom.xml`.

```
<properties>
  <libmemcached.home>${RDMA_MEMCACHED_INSTALL_DIR}/libmemcached
</libmemcached.home>
</properties>
```

2. Run Maven to build the OHB Micro-benchmark for Memcached

```
$ mvn clean package
```

All micro-benchmarks will be installed in `OHB_INSTALL_PATH/memcached/target`. More details on building and running the micro-benchmarks are provided in the `README.memcached.txt`.

This micro-benchmark suite is also available with RDMA-Memcached package in the following directory: `<RDMA-MEMCACHED_INSTALL_PATH>/libmemcached/bin`. A brief description of the benchmark in the following sections.

7.1.1 Memcached Latency Micro-benchmark (`ohb_memlat`)

This micro-benchmark measures latency of a memcached operation for different data sizes. The test can run in three modes:

1. GET - The OHB Get Micro-benchmark measures the average latency of memcached get operation.
2. SET - The OHB Set Micro-benchmark measures the average latency of memcached set operation.
3. MIX - The OHB Mix Micro-benchmark measures the average latency per operation with a get/set mix of 90:10.

In all three micro-benchmarks, the memcached operations are repeated for a fixed number of iterations, for different data sizes (1B to 512KB). The average latency per iteration is reported, ignoring any overheads due to start-up.

The test mode can be selected using runtime arguments `--benchmark=<GET|SET|MIX|ALL>`. Similarly, Memcached server can be specified using the runtime argument `--servers=<SERVER[:PORT]>`. Below is a list of all runtime parameters,

--servers=<SERVER[:PORT]>

List which servers you wish to connect to.

--benchmark=<SET|GET|MIX|ALL>

Specify OHB Micro-benchmark type. Supported micro-benchmarks:

SET - OHB Set Micro-benchmark.

Micro-benchmark for memcached set operations.

GET - OHB Get Micro-benchmark.

Micro-benchmark for memcached get operations.

MIX - OHB Mix Micro-benchmark.

Micro-benchmark for memcached set/get mix.

ALL - Run all three OHB Micro-benchmarks.

--version

Display the version of the application and then exit.

--help

Display help and then exit.

An example of running this micro-benchmark is as follows:

```
$> <LIBMEMCACHED_INSTALL_PATH>/bin/ohb_memlat --servers=storage01:11211
--benchmark=MIX
```

7.1.2 Memcached Hybrid Micro-benchmark (ohb_memhybrid)

This micro-benchmark measures average latency and success rate of memcached gets for different data sizes, for a specified penalty for a miss in the Memcached server. The micro-benchmark first populates Memcached server with more keys than what can probably fit in memory, and accesses these keys at random in one of the two modes:

1. UNIFORM - All keys are selected uniformly at random.
2. NORMAL - Some keys are accessed more frequently than others.

For any of these modes, tests can be run with different spill factors, maximum memory available to Memcached server, and value size of key/value pair. The size of key is fixed at 16 Bytes. Below is a list of all runtime parameters:

--servers=<SERVER[:PORT]>

List which servers you wish to connect to.

--spillfactor=<value>

Spill factor (value greater than 1.0 preferably to simulate over-flow). Default is 1.33.

--maxmemory=<value in MB>

Max Server Memory (in MB). Default is 64 MB.

--valsize=<value in Bytes>

Value Size of Key/Value Pairs (in Bytes). Default is 4KB. Key size fixed to 16 Bytes.

--scanmode=<UNIFORM | NORMAL>

Pattern for memcached_get. Default is UNIFORM. Supported modes:

UNIFORM - All keys are selected at random with equal probability.

NORMAL - Some keys queried more frequently than others (similar to normal distribution).

--misspenalty=<value in ms>

Additional latency (e.g. database access latency) to fetch key/value pair, when it is a miss in memcached (in ms). Default is 1 ms.

An example of running this micro-benchmark where the client populates the Memcached server with 50% more key/values pairs (with value size 4KB) than what can fit into a server running with 64MB of RAM, assuming that the additional latency on a miss at the Memcached layer is 1.5 ms, with data being accessed in a normal pattern is as follows:

```
$$> <LIBMEMCACHED_INSTALL_PATH>/bin/ohb_memhybrid --spillfactor=1.5  
--servers=storage01:11211 --maxmemory=64 --misspenalty=1.5  
--scanmode=NORMAL --valsize=4096
```

7.1.3 Memcached Latency Micro-benchmark for Non-Blocking APIs (ohb_memlat_nb)

This micro-benchmark measures average latency of the newly introduced non-blocking set and get operations, i.e., iset, iget, bset, bget etc. The micro-benchmark issues a batch of first non-blocking set/get requests and monitors their completion using either wait or test APIs. This parameter is referred to as request threshold and behaves as a barrier to ensure progress and completion of all ongoing requests. The micro-benchmark behaves as follows:

1. For micro-benchmark runs with iset or bset, the number of set requests issued is equal to the <max-memory>/<value-size>.
2. For micro-benchmark runs with iget and bget, the Memcached servers are populated with <max-memory>/<value-size> key/value pairs. These pairs can be read using either at random or using a zipf pattern and the number of such requests can be controlled during runtime.

For any of these modes, tests can be run with different aggregated maximum server memory, value size of key/value pair, progress API of choice and the threshold of number of ongoing requests. The size of key is fixed at 16 bytes. Below is a list of all runtime parameters:

--servers=<SERVER[:PORT]>
List which servers you wish to connect to.

--reqtype=<API short name>
iset/bset/iget/bget.

--progresstype=<non-blocking progress API>
wait/test.

--pattern=<pattern for iget/bget>
Pattern for memcached_get/iget/bget. Default is "random". random: All keys are selected at random with equal probability. zipf: Keys selected using zipf distribution.

--maxmemory=<value in MB>
Max Server Memory (in MB). Default is 64 MB.

--numgets=<number of iterations>
Total number of key/value pairs to fetch from Memcached servers.

--valsize=<value in Bytes>
Value size of key/value pairs (in bytes). Default is 4KB. Key size fixed to 16 Bytes.

--reqthresh=<number of requests>
Number of pending non-blocking requests allowed.
Must be smaller than <max-memory>/<value-size>.

--verbose
Display more verbose output for micro-benchmark.

--version
Display the version of the application and then exit.

--help
Display help and then exit.

An example of running this micro-benchmark where the client populates the Memcached server with value size 4 KB and maximum memory of 1 GB, i.e., 256 K key/value pairs using iset API (`memcached_iset`) and wait API (`memcached_wait`) with a request threshold of 32 ongoing requests is as follows:

```
$$> <LIBMEMCACHED_INSTALL_PATH>/bin/ohb_memlat_nb --maxmemory=1024  
--valsize=4096 --servers=storage01:11211 --reqthresh=32  
--reqtype=iset --progresstype=wait
```

An example of running this micro-benchmark where the client reads 1000 key/value pairs from the Memcached server of value size 256 KB and maximum memory of 1 GB, i.e., from 4 K key/value pairs using bget API (`memcached_bget`) and wait API (`memcached_test`) with a request threshold of 16 ongoing requests is as follows:

```
$$> <LIBMEMCACHED_INSTALL_PATH>/bin/ohb_memlat_nb --maxmemory=1024  
--valsize=262144 --servers=storage01:11211 --numgets=1000  
--reqtype=bget --progresstype=test --reqthresh=16
```


7.2 YCSB extension for RDMA-Memcached in OHB

This is an extension to the popular Yahoo! Cloud System Benchmark (YCSB) benchmark available [here](#) for evaluating Memcached. The YCSB Memcached client uses the Spymemcached Java library for running with the TCP/IP-based Memcached servers (<https://github.com/memcached/memcached>). We extend this open-source benchmark to support the RDMA-Libmemcached library-based client through the Java Native Interface (JNI).

7.2.1 Building YCSB with RDMA-Libmemcached Support

Pre-requisites for building the YCSB benchmark suite are listed below:

1. Maven (for building the YCSB-Memcached suite and extension)
2. Java version ≥ 1.7 (JDK installation)
3. RDMA-Memcached software package (RDMA-Libmemcached library)

Steps for building YCSB benchmark suite and JNI layer for MemcachedClient are as follows:

1. Set YCSB home and RDMA-Libmemcached package paths as follows:

```
$$> export YCSB_HOME=<path-to-OHB-install>/memcached/ycsb-  
memcached-ext/
```

```
$$> export LIBMEMCACHED_HOME=<path-to-RDMA-Memcached-install>/  
libmemcached
```

2. Ensure JAVA_HOME is set to the JDK installation path.

```
export JAVA_HOME=<path-to-jdk-install>
```

This is required by the following scripts to build the C code with JNI generated headers: `build.sh` and `clean.sh`

3. Build the YCSB benchmark and JNI with maven

```
mvn clean package
```

NOTE: Both Micro-benchmarks and the YCSB extension can be build together for RDMA-Memcached/RDMA-Libmemcached from `$OHB_INSTALL_PATH/memcached/` by running `'mvn install'`. Ensure that `libmemcached.home` property (step (1) of Section 7.1) and `YCSB_HOME` environment variable (step (1) of Section 7.2) are set before building.

7.2.2 Running the YCSB Benchmark Extension with RDMA-Memcached

Steps for running YCSB benchmark suite with RDMA-Memcached are as follows:

1. Set up a RDMA-Memcached server cluster to benchmark. See Section 5 for more details.
2. Core properties, such as the read-update pattern, record sizes, etc., can be set in the workload configuration file as done in the basic YCSB benchmark. E.g., e.g., `$YCSB_HOME/workloads/workloada`. More details can be found in the [YCSB wiki](#).
3. Set up following RDMA-Memcached specific YCSB workload configurations, e.g., `workloads/workloada`, by adding these two parameters:

- Update the server list

```
memcached.hosts=<comma-separater-list-of-servers>
#e.g., memcached.hosts=server1:port,server2:port,server3:port
```

- Update the protocol to be used

```
memcached.protocol=RDMA #set protocol to RDMA
```

- Java Direct Byte Buffer is used to pass data between RDMA-Libmemcached and YCSB. This represents number of individual records (or buffers) to allocate.

```
memcached.native.nbuffer=<number>
# e.g., memcached.native.nbuffer=4
```

- Buffer slots are allocated with size: `MAX_KEY_SIZE + PADDING + (field_count * field_length)`. This parameter can be used to change `MAX_KEY_SIZE`. Default is 256 bytes; this is the Memcached protocol maximum.

```
memcached.key.maxsize=<number>
# e.g., memcached.key.maxsize=128
```

Detailed configuration with example is provided in `$YCSB_HOME/memcached/workloads/sample_workload`.

4. Run YCSB workload with `LD_LIBRARY_PATH` to point to RDMA-Libmemcached library and native JNI-based library build for `RdmaMemcachedClient` class (built using above steps):

- Change diretory to YCSB extension home directory:

```
cd $YCSB_HOME
```

- Load Memcached servers with records:

```
LD_LIBRARY_PATH=${LIBMEMCACHED_HOME}/lib:${YCSB_HOME}/  
memcached/target/native_libs:${LD_LIBRARY_PATH} bin/ycsb.sh  
load basic -P workloads/sample_workload
```

- Run the YCSB workload:

```
LD_LIBRARY_PATH=${LIBMEMCACHED_HOME}/lib:${YCSB_HOME}/  
memcached/target/native_libs:${LD_LIBRARY_PATH} bin/ycsb.sh  
run basic -P workloads/sample_workload
```

Running the ‘ycsb’ command without any argument will print the usage.

NOTE: Using Spymemcached for TCP/IP-based Memcached servers: The default Memcached java client can be enabled to run with default Memcached servers by setting the protocol in the workload configuration:

```
memcached.protocol=TEXT #or BINARY
```

This will revert the benchmark back to using Spymemcached.

8 Troubleshooting with RDMA-Memcached

If you are experiencing any problems with RDMA-Memcached package, please feel free to contact us by sending an email to rdma-memcached-discuss@cse.ohio-state.edu.