

A REPORT ON A KADEMLIA IMPLEMENTATION

---

# KADEMLIA

---

April 5th, 2019

Christopher Mykota-Reid  
Rowan MacLachlan

*CMPT 434: Computer Networks*

Derek Eager

# Contents

0.1	Introduction . . . . .	3
0.2	Objectives . . . . .	3
0.3	Implementation . . . . .	4
0.3.1	Use . . . . .	4
0.3.2	High-Level Overview and Tooling . . . . .	5
0.3.3	Kademlia Remote Procedure Calls . . . . .	8
0.4	Measurements and Criteria . . . . .	10
0.5	Technical Difficulties . . . . .	11

## 0.1 Introduction

DHTs are similar in principle to a regular hash table in that they provide the structure for fast retrieval of an object by referencing the location of the object's hash value. However, a distributed hash table is stored across nodes in a network, and this introduces a variety of complications involving node-lookup (finding the location of the stored data from its hash in the network), ensuring data redundancy (what happens if a node goes offline?), and data retrieval. Although this technology is not *new* by any standard, its continued use and development in various fields makes it a suitable candidate for an implementation project.

Implementing the Kademlia protocol provided an excellent learning opportunity with respect to various course materials:

- practical application-level network-related programming,
- the use of some *\*NIX* network tools (*ifconfig*, *ip*, *netstat*, etc.),
- the application of some course theory (implementing application-layer communication across different networks over UDP implies some difficulties with port-forwarding at the router — this is what NAT would help with,)
- a much deeper understanding of the Kademlia protocol, and
- an understanding of some of the issues overlay networks (DHTs in particular) must deal with.

## 0.2 Objectives

As referenced above, the problems presented by a DHT implementation are myriad and often difficult to solve, but by no means impossible. By implementing an incremental plan to gradually roll out more complex features, taking full advantage of existing libraries, and using a high-level and expressive programming language, we managed to implement a DHT. While this DHT does not implement every particular directive of the Kademlia white paper exactly (nor is every successfully implemented feature necessarily bug-free), we managed to create an application that could be easily run and could, without issue, support a small network of up to 8 devices.

On one level, the implementation team is very happy with this result, considering the considerable hours invested in research, implementation, and testing. Although we failed to produce a working proto-type with enough left-over time to collect, interpret, and display data from large-scale tests (a hundred or more nodes), we feel as though the preliminary testing and class demo show at least some acceptable level of achievement. While this may be disappointing, it must be noted that many excellent Kademlia implementations already exist, most of which feature more finely-tuned, bugless, and efficient code. In this sense, simulation data would only have cemented an already clear understanding of our implementation's shortcomings. Implementation decisions were almost

always made by giving far more weight to deadlines than to writing performant code. After all, premature optimization is the root of all evil.

More importantly, we exceeded our personal academic goals — with respect to interest, learning, and technical achievement, the project was a complete success.

## 0.3 Implementation

### 0.3.1 Use

Using the program is relatively simple. The project folder, once downloaded, either provided by the students or downloaded from the GIT repository here, has a simple structure. The source code is in the `kademlia` folder. So long as the user has a python3.7 (or greater) environment, everything should work.

1. Configure your python environment so that the `python3` command points to an installation of python3.7 or greater.
2. `git clone https://github.com/rowan-maclachlan/cmpt-434-proj`
3. `cd cmpt-434-proj`
4. `make init`. This will install the application dependencies.
5. `python3 kad.py`. This will launch a single Kademlia node. To see application use, run the program without arguments.
6. Use command `set`, `get`, `ping`, and `inspect` to interact with the DHT.

In addition to this relatively simple use, the repository also contains a script to launch an arbitrary number of nodes to run in a non-interactive way. More details exist within the script itself, but this simulation script is operated as such: `./simulation.sh 10 127.0.0.1 2000`. This launches 10 nodes that bootstrap off of an already existing node located at `127.0.0.1:2000`, but these nodes don't accept user input. Instead, they can be used to test the correctness and robustness of the system.

Please view the README for further details.

## 0.3.2 High-Level Overview and Tooling

### 0.3.2.1 Networking

The Kademlia paper implies that the protocol is implemented with datagrams[1]. Typically, clients operating behind a NAT (Network Address Translator) do not need to worry about querying servers outside the NAT, as the Network Address Translator maintains a table correlating incoming messages to hosts on its network. However, on peer-to-peer systems, clients must also act as servers, and accept requests from other nodes on the overlay network despite being located behind a NAT. This requires the use of a NAT traversal technique such as Hole-Punching or SOCKS (Socket Secure) to overcome the difficulties inherent in routing requests from LANs over the internet. This makes application use on LANs behind a router address difficult without port-forwarding and configuration overhead. Because of this, our application was tested only within a LAN.

### 0.3.2.2 RPC library

Our primary investigation showed other Kademlia implementations using the core python library `asyncio` to manage callback procedures for asynchronous code execution. While the semantics of asynchronous RPC might make implementation *more* difficult than for a simple synchronous I/O implementation (the python `socket` module directly, for example), implementing asynchronous UDP RCP with `asyncio` is very simple. In fact, we used a 3rd party library called `rpcudp` which overlays functionality for remote procedure calls onto a datagram communication protocol - all asynchronously![2]

As described on the `asyncio` Read-The-Docs page: "asyncio is often a perfect fit for IO-bound and high-level structured network code." This is exactly the Kademlia use-case.[3]

### 0.3.2.3 Hashing

Kademlia specifies the use of the 160-bit SHA-1 hash for data, and 160-bit node IDs and 160 `k-buckets` (0.3.2.5), to accompany that. This parameter of ID and hash length, however, is not mandated. The python library `hashlib` provides an array of hashing functions, cryptographic and otherwise, including the SHA-1 hash. This library can be used behind a small wrapper, and hash/ID size can be trimmed to the length specified in the simulation parameters. This value could initially be set quite small to increase ease-of-testing and provide more comprehensible and tractable output during development stages. It can easily be changed later.

```
def hash_function(data):  
    """  
    Hash the data to a byte array of length  $p.params[B] / 8$   
  
    Parameters
```

```
data : binary data
```

```
Return
```

```
int : The digest of the input data, trimmed to the ID space.  
"""
```

```
return int(hashlib.sha1(data).hexdigest(), 16) & get_mask()
```

As shown here, the returned digest is ANDed with a mask — this mask is the length of the system-wide hash/ID bit length value — called ‘B’ in the Kademlia paper. This mask then truncates the SHA-1 hash of the data to a value between 0 and  $2^b - 1$ . This allows the value of ‘B’ to be effortlessly changed for more practical use of the network or for larger-scale simulations.

#### 0.3.2.4 Routing Table Design

Study of different Kademlia implementations[4][5][6][7] reveal that most implementations implement their routing table in a way very similar to that described in the Kademlia white paper[1]. This involves creating nodes with only a single `k-bucket` 0.3.2.5. As the `k-bucket` reaches capacity, it is then split into two and its contents distributed among the two resultant buckets according to their proximity to the node that owns the routing table. In this sense, a routing table may not actually have as many `k-bucket`s in memory as there are bits in the keyspace. Naturally, this is far more space efficient than the implementation we opted for.

As the maximum number of `k-bucket`s on any particular node is simply the length of the ID space, they can all be created at once, as in the initialization code for a routing table 1:

**Listing 1:** RoutingTable()

```
self.buckets = [ KBucket(k) for _ in range(b) ]
```

Instead of attempting to insert nodes into a `k-bucket` which may or may not be full and which may result in further operations on the data structure, nodes can be inserted directly into the `k-bucket` corresponding to their distance from the routing table’s owner, as shown in 2:

**Listing 2:** RoutingTable.add()

```
return self.get_bucket(contact.getId()).add(contact)
```

where the method `RoutingTable.get_bucket(id)` calculates the correct bucket index by referencing the most significant bit of the distance between the ID of the routing table’s owner and the ID of the contact we are trying to add to the routing table, like so:

**Listing 3:** RoutingTable.get\_bucket()

```
distance = self.id ^ id
index = 0
while distance > 1:
    # Count the index of the largest bit in the distance
    # The distance will be zero after 'bit' many shifts.
    distance = distance >> 1
    index += 1
return self.buckets[index]
```

Naturally, `k-bucket`s can be retrieved in a similar manner.

Although this approach is less space efficient than creating `k-bucket`s only as needed, the overhead is fixed and relatively small. Ultimately, this decision was made because it was believed to be simpler.

### 0.3.2.5 Bucket Design

The `k-bucket` is a container that holds contact information for other nodes on the Kademlia network. Each `k-bucket` is created with a maximum capacity of `k`, a system-wide parameter that controls for such things as `k-bucket` size. In our implementation, it is simply a light wrapper for a list into which contacts are stored.

As described by the white paper, `k-bucket`s, once full, need to make evaluations about whether or not new contacts are inserted into the `k-bucket` after it has reached capacity. When a new node is discovered but its `k-bucket` is full, it is only added if the node at the head of the list is unresponsive. Otherwise, the new node is simply discarded, and the responsive node is moved from the head to the tail of the list.

There are a variety of reasons for this behaviour. First, the Kademlia authors illustrate that nodes which have existed on the network for a long time are statistically more likely to continue being active on the network than new nodes. Therefore, the integrity of the network as a whole is improved by prioritizing these old nodes. Secondly, this behaviour makes the network automatically resistant to Denial-of-Service attacks: the network cannot be brought down by flooding it with new nodes, because the old network will not replace legitimate nodes that already exist in its routing table unless they become unresponsive.

Unfortunately, our implementation does not track node responsiveness, as the minimum-viable-product was simpler if new nodes are simply added (because they are known to be active) and old nodes are kicked out of the routing table.

**Listing 4:** KBucket.add()

```
if contact in self.contacts:
```

```

        # If this contact already appears in the list , move it to the back
        # of the list
        self.contacts.remove(contact)
        self.contacts.append(contact)
    elif not self.full():
        # If the contact doesn't appear in the list , append it.
        self.contacts.append(contact)
    else:
        # If the list is full , remove the oldest contact before appending
        # the new one.
        self.contacts.popleft()
        self.contacts.append(contact)
    return True

```

Because most operations on the `k-bucket` involve adding and removing nodes from the beginning and end of the list, a special-use data structure was employed to optimize the efficiency of these operations. Python provides a structure called a ‘deque’ (pronounced ‘deck’) that implements constant-time additions and removals from the head and tail of the list.

### 0.3.3 Kademlia Remote Procedure Calls

#### 0.3.3.1 Store

The store operation is used to save a value onto the Kademlia network. Our application accepts a `String`-type key and a `String`-type value, as such: `set my-key my-value`. The key digest is then taken, which maps the key to a value in the node-space range (between 0 and  $2^b - 1$ ). The node then attempts to find the closest node on the network to the key digest via the `find_nodes` 0.3.3.5 algorithm. It then issues a store request to  $k$  of the nodes closest to the value which it found in its search. A store request message consists of the sender’s ID, the digest, and the value. A node that receives a store request assumes it is responsible for the key and stores the value without exception. Because the sender ID is sent along with the message, the receiving node can take steps to properly handle the new contact via the `handle_node` procedure (see 5).

#### 0.3.3.2 Find Value

This operation is used to find a value on the network. Implementation depends on the iterative Node Find algorithm defined below 0.3.3.5. When this command is received by a node, it returns to the sender either the value requested (if the receiver has the value in its data store) or a list of the  $k$  closest nodes in its routing table to the value being sought.



### 0.3.3.3 Find Node

This operation is very similar to the Find Value 0.3.3.2 operation. The difference is that the node issuing the request does not need to anticipate receiving a value *or* a list of contact — only contacts are returned by the receiving node.

### 0.3.3.4 Ping

This operation is used to check for liveness on the network. Specifically, it is used by nodes to determine whether a new contact should replace an old contact in its routing table. If the old contact is unresponsive (it does not respond to the `ping` operation), it will be removed from the routing table in favour of the new contact, which is known to be live. This is described in more detail in the implementation section 0.3.2.5 on `KBucket` design.

### 0.3.3.5 Finding Nodes on the Network

Searching on the Kademlia Network is done using an iterative, parallel algorithm. There are 2 parameters that control the search:  $\alpha$  and  $k$  (the same  $k$  referred to in k-buckets).  $\alpha$  controls the number of parallel queries active at once and  $k$  controls the number of nodes returned by each query of a node as well as the maximum number of nodes queried before the search terminates. The nodes queried are taken from a list of nodes, ordered by their distance to the target, called the shortlist. The algorithm works by selecting the  $\alpha$  closest nodes to the target from the shortlist and querying them for nodes. These queried nodes respond with their  $k$  closest contacts to the target (retrieved from their routing table), which could contain the target itself. The contacts are then put into the shortlist to be selected for future queries. If no node is found closer to the closest node found so far then the  $k$  closest nodes not yet queried are queried, resulting in the termination of the search after the responses are processed.

This process is repeated until one of three conditions is met:  $k$  nodes have been queried, the shortlist is empty, or we have found the target we were searching for. Upon completion of the search, the node that called the search is returned either:

- the  $k$  closest nodes to the target found, or,
- a value, if the search was a value search.

The other search used in Kademlia is a Value Search. This search functions almost identically to the Node Search — differing only if the target is found.

- If the target is found in a Node Search, the  $k$  closest contacts to the target are returned.

- In a Value Search, only the value associated with the target is returned. A store 0.3.3.1 is then issued to the closest node found in the search that did not have the target data.

If a Value Search fails (probably because the Kademlia network does not contain a mapping for the key), then no store is performed and the  $k$  closest nodes to the target are returned, just as in the Node Search.

## 0.4 Measurements and Criteria

As referenced previously (0.2), we failed to collect any meaningful simulation data, and so network performance was not measured. The only quantitative measurement of project success we have is the number of active nodes we could launch on the same network. Preliminary testing revealed we could successfully launch as many as 10 nodes on the same network, and observe correct performance of node discovery, set, and get operations. This was achieved through a small script to launch multiple nodes at once, called `simulation.sh`.

As described in 0.3.1, using this script allows the user to make an evaluation of the system behaviour, robustness, and correctness. Although no meaningful performance data are collected from this simulation, the infrastructure is in place. Here is the result of running the `inspect` command on a node used to bootstrapping 40 others:

```
'set <key (str)> <value (str)>' to store data
'get <value (str)>' to retrieve data
'inspect' to view this node's state
'quit' to leave

Attempting to run inspect...
Data for this node: {}
Routing table for (38038,0.0.0.0,2000)
10: deque([(37067,127.0.0.1,2011)], maxlen=4)
11: deque([(38917,127.0.0.1,2014), (40255,127.0.0.1,2028), \
(40086,127.0.0.1,2038)], maxlen=4)
12: deque([(35664,127.0.0.1,2019)], maxlen=4)
13: deque([(41656,127.0.0.1,2026), (43503,127.0.0.1,2033), \
(45280,127.0.0.1,2035), (44437,127.0.0.1,2036)], maxlen=4)
14: deque([(55477,127.0.0.1,2037), (61319,127.0.0.1,2040), \
(54538,127.0.0.1,2044), (49313,127.0.0.1,2049)], maxlen=4)
15: deque([(27086,127.0.0.1,2046), (8111,127.0.0.1,2047), \
(32276,127.0.0.1,2048), (14986,127.0.0.1,2050)], maxlen=4)
```

As we would expect, higher-index `k-buckets` are fully populated, while the fraction of nodes which reside within the close neighbourhood of node 38038 is small. As can be derived from the high-number `k-bucket` index, the id-space of this simulation is  $b = 16$ .

## 0.5 Technical Difficulties

Here, we list some of the technical considerations and their resolutions:

- How many nodes store the data of a single object? What degree of redundancy can be achieved and what overhead does this add?

The Kademlia paper invokes this concern, but does not provide any hard-numbers. According to the authors, this number must be high enough that even on very difficult networks, no more than this number of nodes in the same key-space go offline within an hour of each other. The authors provide a value of 20, but we tested with a far smaller value so that correct behaviour could be more easily observed.

- Is an entire object stored on each node, or do we distribute a single object in incomplete parts through the network?

Our implementation accepted only `String` types as keys and values. Regardless of the message type (node list, node id, key, or value) the entire message is always transmitted in a single UDP packet. This limits message size to a theoretical value of 65,507 bytes minus the overhead incurred by the `rpcudp` library, but our implementation did not account for message sizes in any way. Presumably, the user could crash their node by trying to store a massive value.

- How large is the address space: how large are the hash keys? 128 bits? 160 bits? We may not need to have a large address space for our implementation, but just what are the consequences of this choice?

As described in 0.3.2.3, key-space size was easily adjustable because of system design. However, the decision was made during implementation to keep the value small. This made it easy to manually verify correct network behaviour by calculating distances by hand, as well as put more stress onto some system internals such as 0.3.2.5 and 0.3.2.4 which behave differently depending on the density of nodes around the keyspace.

- Do we hash the object's identifier or the object data? If we hash the identifier, how do we enforce unique identifiers?

Unique identifiers (keys) are not enforced. If some aspect of uniqueness was taken from the node attempting the store, their ID or IP/port could be used to create a unique hash. However, This would make it difficult for other nodes to request that data (as the hash was created with some other node's identifying information.) A better solution would be to request the data mapped to by the hash-key before attempting the store, and only allow the store if the data was not found. While this solution is not technically difficult, it was deemed unimportant in relation to other features.

- Any hash function will have collisions. How do we manage collisions? Do we chain hash contents at the site of storage or do we disallow colliding files?

Support for chaining hash-collisions was not implemented. Data storage was managed simply by Python's built-in dictionary type. If there was a collision (highly possible on a test-network with a smaller than usual key-space) the previous value would be overwritten at the node performing the store. This could be avoided by the use of a more powerful container type, but it would introduce other questions and corresponding technical issues.

When requesting a value from a key which maps to a node containing chained occurrences of that key, which value should be returned? There's no way for the node receiving the request to know which value the node making the request is after, so should it return both values? This would throw a wrench in the current implementation of the `get` command, which does not expect more than a single value response.

- How does the *CAP* theorem (or *Brewer* Theorem) inform our implementation? What does it mean for it to be impossible for a distributed system to provide all of Consistency, Availability, and Partition Tolerance?

As described by the theorem, providing both consistency and availability on well-managed internal networks (our test-case) is generally not an issue. Even so, the Kademlia protocol does an excellent job at specifying the responsibilities of nodes to keep their routing tables up-to-date and ensure that the key-values they are responsible for exist with a high degree of redundancy. In the event of highly-congested networks or networks with partitions, the Kademlia protocol (with some small changes to our existing implementation) will naturally route around high-latency links. This is achieved by the continuous re-ordering of active and responsive nodes in the routing table, and the option of storing contacts within k-buckets ordered on round-trip-time, which can be easily measured.

- How did we manage the issues presented by key-space remapping?

Is it possible for us to avoid the issue of key-space re-mapping (changing the node location of data) when adding or removing nodes from the network?

Thankfully, this is handled very nicely by the Kademlia protocol. The first issue it resolves is re-mapping existing data to a new node, and the second issue is caused by nodes leaving the network. Imagine a node  $n_1$  joining the Kademlia network. It does this by pinging another node on the network that it knows about,  $n_2$ . When  $n_2$  hears about  $n_1$  (and  $n_1$ 's ID),  $n_2$  issues `store` requests to  $n_1$  for every piece of data  $n_2$  has in its routing table for which the following conditions are met:

- $n_1$ 's ID is closer to the data (by XOR) than the furthest among our k-many closest already known contacts to the data.
- We are closest to the data of any of our already known contacts.

This ensures that new nodes are holding data they are responsible for, and it also minimizes the number of extraneous stores being sent over the network. This is perhaps difficult to understand but section 2.5 of the Kademlia paper describes it well. Or, consider the code below:

**Listing 5:** Protocol.handle\_node()

```
def handle_node(self, contact):
    ...
    # See Kademlia paper section 2.5 on how to incorporate new nodes
    # We may have to store all values we have which are closer to
    # the
    # new node than they are to us.
    for key, value in self.data.items():
        # find neighbours close to the key value
        nearest_contacts = self.table.find_nearest_neighbours(key)
        # If there are fewer than k neighbours, store the key-value
        # to the
        # new node
        if len(nearest_contacts) < self.table.k:
            log.debug(f"Few contacts, storing data to new contact..."
                    ")
            # schedule the task in the event loop, continue to next
            # data
            asyncio.create_task(self.try_store_value(contact, key,
                value))
            continue
        # If there are k neighbours, only store the key-value if the
        # new
        # node is closer to the key than our neighbour furthest from
        # the
        # key, and if we are closer to the new node than any of our
        # neighbours.
        nearest_contact = nearest_contacts[0]
        # are we nearer to this key than nearest_contact is?
        were_nearest = \
            self.this_node.distance(contact.getId()) < \
            nearest_contact.distance(contact.getId())
        furthest_contact = nearest_contacts[-1]
        # Is contact closer to the key than the furthest contact?
        close_enough_to_store = \
            contact.distance(key) < furthest_contact.distance(
                key)
        if close_enough_to_store and were_nearest:
            log.debug(f"Storing {data} to new contact...")
            asyncio.create_task(self.try_store_value(contact, key,
                value))

    self.table.add(contact)

    return True
```

The second issue of nodes leaving the network is handled by there being a degree of redundancy for data. Data are replicated by a factor of  $k$  on the network, and so a single node leaving has no practical effect.

Please see the bibliography for additional reading.

# Bibliography

- [1] Petar Maymounkov and David Mazieres. “Kademlia: A peer-to-peer information system based on the xor metric”. In: *International Workshop on Peer-to-Peer Systems*. Springer. 2002, pp. 53–65.
- [2] Brian Muller. *RPC over UDP Python module*. URL: <https://github.com/bmuller/rpcudp> (visited on 04/02/2019).
- [3] *Python asyncio module documentation*. URL: <https://docs.python.org/3/library/asyncio.html> (visited on 04/02/2019).
- [4] Brian Muller. *Kademlia Implementation*. URL: <https://github.com/bmuller/kademlia> (visited on 04/02/2019).
- [5] Isaac Zafuta. *Kademlia Implementation*. URL: <https://github.com/isaaczafuta/pydht> (visited on 04/02/2019).
- [6] flosch. *Kademlia Implementation*. URL: <https://github.com/flosch/libdht> (visited on 04/02/2019).
- [7] cjgu. *Kademlia Implementation*. URL: <https://github.com/cjgu/asyncio-kademlia> (visited on 04/02/2019).
- [8] Stefan Saroiu, P Krishna Gummadi, and Steven D Gribble. “Measurement study of peer-to-peer file sharing systems”. In: *Multimedia computing and networking 2002*. Vol. 4673. International Society for Optics and Photonics. 2001, pp. 156–171.
- [9] Ion Stoica et al. “Chord: a scalable peer-to-peer lookup protocol for internet applications”. In: *IEEE/ACM Transactions on Networking (TON)* 11.1 (2003), pp. 17–32.
- [10] Antony Rowstron and Peter Druschel. “Pastry: Scalable, decentralized object location, and routing for large-scale peer-to-peer systems”. In: *IFIP/ACM International Conference on Distributed Systems Platforms and Open Distributed Processing*. Springer. 2001, pp. 329–350.
- [11] Ben Y Zhao et al. “Tapestry: A resilient global-scale overlay for service deployment”. In: *IEEE Journal on selected areas in communications* 22.1 (2004), pp. 41–53.
- [12] *BitTorrent*. URL: <https://en.wikipedia.org/wiki/BitTorrent> (visited on 03/21/2019).
- [13] *Kademlia Specification Guide*. URL: <http://xlattice.sourceforge.net/components/protocol/kademlia/specs.html> (visited on 04/02/2019).