

# Distributed Key-Value Storage Algorithm

## CS181E — Distributed Systems

### Assignment 5

Alejandro Frias

Ravi Kumar

David Scott

April 16, 2014

## Algorithm Setup

In this problem, we have a setup of Erlang nodes arranged in a ring. Communications may occur across “chords” of the ring according to specifications laid out in the assignment. The important part of the communication restrictions for the purposes of our implementation is that a non-storage process on a given node is guaranteed to be able to send messages to a non-storage process on the node after it in the ring.

## Types of Processes

Our system involves two types of processes: `storage_processes` and `store_handlers`. We have the `storage_processes` responsible for the primary data storage in the system (both storing and retrieving), and the `store_handlers` responsible for meta-storage functions including the backing up of data and the handling of snapshot-related functionality. On a given Erlang node, there will always be exactly one handler running, and some number of `storage_processes`. In our system, the data on a `storage_process` is backed up by the handler on the node “in front” of the processes node.

## Storage Processes

`storage_processes` in our system are responsible for receiving all messages from the outside world. When a `storage_process` gets a store message, they forward it to the closest `storage_process` to the destination. When the correct `storage_process` gets the store request, it stores the data and notifies its `store_handler`. When an `storage_process` gets a retrieve message, it forwards it along until it gets to the correct `storage_process`, which reports to the outside world. Finally, `storage_processes` can easily handle `node_list` messages, as they have access to the global registry. All other messages received from the outside world are simply forwarded to the `storage_process`’s handler. A `storage_process` has only one possible state, that of waiting for messages. For a more detailed de-

scription of what messages a `storage_process` sends and receives and what actions it takes upon receiving these messages, refer to the message description table.

In our system, a `storage_process` stores the following information:

**m** The size of the system, so that it knows what other processes it can message.

**myID** The ID of the `storage_process`.

**myDict** A dictionary of data stored by the process.

**myHandlerID** The ID of the node's handler, so that it can message it.

## Handler

Handlers are responsible for everything else in the system. They handle all requests related to system snapshots, manage when `storage_processes` should be started or stopped due to rebalancing when nodes join and leave, and keep backups of data. Specifically, the `storage_process` for a given node backs up all data for the node behind it. If the node behind it dies, it will use this backup data to start all `storage_processes` that died, ensuring that no data is lost. For a more detailed description of what messages a `store_handler` sends and receives and what actions it takes upon receiving these messages, refer to the message description table.

In our system, a `store_handler` stores the following information:

**m** The size of the system.

**myID** The ID of the handler process.

**nextNodeID** The ID of the next handler in the system. The handler uses this to send messages around the ring and to send data to the next node to be backed up.

**prevNodeID** The ID of the previous handler in the system, which this node is monitoring for failure.

**myBackup** A backup of the data held by all `storage_processes` on the previous node.

**minKey** The minimum key in the backup data.

**maxKey** The maximum key in the backup data.

**myBackupSize** The size of the backup data.

**myInProgressRefs** A list of snapshot computations we have started by sending them around the ring. If we get a snapshot message with a ref in this list, we know the computation is finished and can send the result to the outside world.

**myMonitoredNode** The node object we are monitoring for failure.

## Algorithm Description & Correctness

The setup: Each node has a backup of the node preceding it, held by a non-storage process for that node that we call the **store\_handler**, since it also handles much of the inter-node communication. Each node is listening to the preceding node for crashes. In this way each node has one node that it is responsible for and one that is responsible for it.

When a node enters the system, it first finds the location in the system it should be at. To do so, it finds the node that has the most **storage\_processes** and joins such that the new node takes over half of these processes. The new node grabs the entire backup data of the node in front of it (since this is backing up the data that the new node's **storage\_processes** will be responsible for); it also tells next node to delete the first half of its backup data since the new node will now be backing it up. The new node uses the first half of that data (that is, data corresponding to processes behind it) to create its backup and the other half to start the **storage\_processes** it will be responsible for. Additionally, it messages the node behind it to stop running the relevant processes so the new node can run them. In this way, a joining node only needs to talk to two other nodes, the ones that it will be adjacent to. Thus, our system maintains its functionality and correctness when new node's join, as no information is lost and everything stays backed up.

When a node crashes or leaves, the next node finds out, as the next node is the only one listening to it. Since it has the entire backup data of that node, it can immediately startup all of the processes that just died. So it immediately takes over for the dead node, changing its own node number to the one of the deceased. Then, since the data on the nodes that died aren't being backed up, we send a copy of this data to the next node to be backed up. Then, since we lost the backup data that was on the dead node, it sends a request around the ring to the node just previous to it to get all the data on that node's processes. This takes a maximum of  $m$  messages, where there are  $2^m$  processes. In this way we can get the system back up and running immediately and then, while still accepting requests from the outside world, start rebuilding the back up data. When a node dies, the node in front of it takes over for it. This method ensures that our system is correct when nodes die, as the node in front of it has the full backup data for the node that died. Thus, our system works correctly even when nodes die, and doesn't lose any data.

To ensure redundancy, a **storage\_process** never communicates back to the outside world directly when a store request comes in. Instead, after the store request is forwarded to the correct **storage\_process**, that process stores the new key-value pair and tells its handler process that it did so. The handler then sends this to the next node to be backed up. That next node will then notify the outside world that the value has been stored after it stores the backup. This ensures every store request has been backed up before letting the outside world know that the store is completed. Also, very few messages need to be passed. It takes a maximum of  $m$  messages to get the store request to the right **storage\_process** and then a couple more after that to back the data up and respond to the outside world, since the next node is guaranteed to be visible to the node that has that **storage\_process**. This is part of why we chose to have nodes store the back up of the previous node.

Messages about the system, like first-key and last-key, are forwarded directly to the handler process to take care of. The process will start a message that will go around the ring of nodes and will store that node's contribution to calculating the requested computation. Since each handler has a backup of the previous node, it doesn't need to communicate directly with each of the processes, but instead just sees if the first-key that has been found is better or worse than it's first key, or similar comparison for the other key requests. Once a handler gets its message back, it can tell the outside person the result of the snapshot. The runtime of these system processes scales with the number of nodes. They are guaranteed to work if the a node doesn't crash mid computation. Each node's computation will be constant since it can keep track of the first key, the last key, and the number of keys in it's back up as store requests come in. And since a completely stored value is only considered complete once the handler has backed up one up, we're guaranteed to be accurate.

Retrieve requests are the only messages that don't involve a handler at all. They simply forward the message along the closest chord or if they have the value (or should have the value), return the result to the outside world.

Overall, our system's correctness argument is fairly basic. The system works correctly when the system remains static, as store requests work and ensure that the data is both in the correct storage node and backed up before they finish. When nodes join, we ensure that all data remains backed up, so that the state of our systems stays correct. Finally, when nodes die, other nodes immediately step in and replace the dead `storage_processes` using the backup data, and quickly work to replace the backup that was lost, so that all data in the system is backed up. Thus, our distributed hashtable algorithm is correct even in the face of node's joining and leaving.

## Testing

We tested our code in a variety of ways, largely using the test scripts contained in `t.erl`. First, we started a single node with a variety of values for `M`. Then, we did simple tests to ensure that everything was initialized properly, as shown in `test_store_basic`. Then, we insert data into the system, and ensure that everything is correct, such as the number of keys, min and max keys, and that retrieving all the data worked. Then, we tested that pushing new data to keys already in the system would return the correct old values.

After verifying that the basic capability worked, we started testing multiple nodes. We added nodes one by one, repeating the above tests each time and ensuring that everything still worked. Then, we killed nodes arbitrarily until we had again reached one node, checking for correctness each time. Finally, we added and removed nodes again arbitrarily, with more than one.

We also did a series of the above tests with  $m = 2$ , so that we could test a system with as many nodes as storage processes. We basically performed the same tests as above, additionally testing from 4 nodes downwards. We also tested with  $m = 1$ , which caught a bunch of stupid bugs.

# Messages

All of the messages of our system — and the bulk of the algorithm — are described in the very long Table 1. All entries in the right-hand column were intended to be presented in `typewriter` font but L<sup>A</sup>T<sub>E</sub>X was being uncooperative.

Table 1: Description of messages in the system. In this table, we use the initialism “OW” to stand for Outside World: i.e., an arbitrary computer outside of our system.

Message Description	Erlang Pattern
Store Request. Sent by either OW or a <code>storage_process</code> to a <code>storage_process</code> to store <code>value</code> for <code>key</code> . These are only ever received by <code>storage_processes</code> in the network. When a <code>storage_process</code> receives such a message, it first checks if the hashed key is equal to its own ID. If it is, then it stores the value for the specified key, and sends a <code>backup_store</code> message to its <code>store_handler</code> . If not, it forwards it to the closest <code>storage_process</code> to the destination (closest here meaning nearest process <i>before</i> the destination, since processes can only send messages forward in the ring).	{Pid, Ref, store, Key, Value}
Stored Confirmation. Sent by an <code>store_handler</code> to OW after the corresponding request has been stored in the proper <code>storage_process</code> and the data has been backed up in the <code>store_handler</code> sending the message. <code>OldValue</code> is <code>no_value</code> iff this is the first time something has been stored for the given <code>Key</code> .	{Ref, stored, Old-Value}
Retrieve Request. Sent by OW and <code>storage_processes</code> ; received by <code>storage_processes</code> . When a <code>storage_process</code> receives such a message, it checks if the hash of the key is equal to its own process ID. If it is, the <code>storage_process</code> has the value for that key, and replies with a retrieve response. If not, it forwards it to the closest <code>storage_process</code> to the destination (closest as defined above).	{Pid, Ref, retrieve, Key}
Retrieve Response. Sent by <code>storage_processes</code> to the OW. After a <code>storage_process</code> receives a retrieve request meant for it, it looks up the relevant value and reports it to the requesting process in the OW. <code>Value</code> is <code>no_value</code> iff there is no known value corresponding to the given key.	{Ref, retrieved, Value}
First Key Request. Sent by the OW to <code>storage_processes</code> , and by <code>storage_processes</code> to <code>store_handlers</code> . When a <code>storage_process</code> receives this, it will forward the message up to its <code>store_handler</code> . When such a message is received by a <code>store_handler</code> , it will start a first key computation by adding the ref to its list of ongoing computations and sending a First Key Computation message to the <code>store_handler</code> of the next node in the ring.	{Pid, Ref, first_key}

This table continues on the next page...

Table 1: Description of messages in the system. In this table, we use the initialism “OW” to stand for Outside World: i.e., an arbitrary computer outside of our system.

Message Description	Erlang Pattern
Last Key Request. Sent by the OW to <code>storage_processes</code> , and by <code>storage_processes</code> to <code>store_handlers</code> . When a <code>storage_process</code> receives this, it will forward the message up to its <code>store_handler</code> . When such a message is received by a <code>store_handler</code> , it will start a last key computation by adding the ref to its list of ongoing computations and sending a Last Key Computation message to the <code>store_handler</code> of the next node in the ring.	{Pid, Ref, last_key}
Num Keys Request. Sent by the OW to <code>storage_processes</code> , and by <code>storage_processes</code> to <code>store_handlers</code> . When a <code>storage_process</code> receives this, it will forward the message up to its <code>store_handler</code> . When such a message is received by a <code>store_handler</code> , it will start a num keys computation by adding the ref to its list of ongoing computations and sending a Num Keys Computation message to the <code>store_handler</code> of the next node in the ring.	{Pid, Ref, num_keys}
Messages About Keys. Sent and received by <code>store_handler</code> . If Ref is in the list of the receiver’s in-progress computations, the computation is over and the received message contains the result. The receiver will then send the result <code>ComputationSoFar</code> back to the OW. Otherwise, it will update <code>ComputationSoFar</code> with its relevant value and forward the message to the next node’s <code>store_handler</code> .	{Pid, Ref, *_key(s?), ComputationSoFar}
	Where that the atom is a regular expression
Node List Request. Sent by the OW to <code>storage_processes</code> , and by <code>storage_processes</code> to <code>store_handlers</code> . When a <code>storage_process</code> receives this, it will forward the message to its <code>store_handler</code> . When such a message is received by a <code>store_handler</code> , it will query the global registry for the list of nodes, then filter and sort those results to include only the ID numbers of <code>handler</code> nodes. That result is then returned to the requester.	{Pid, Ref, node_list}
Request Result. Sent to the OW by a <code>store_handler</code> . This reports the result of a First Key, Last Key, Num Keys, or Node List Request to the original requester after the <code>store_handlers</code> have finished computing the result.	{Ref, result, Result}
Failure Notification. Sent to the OW by <code>store_handlers</code> or <code>storage_processes</code> to notify the OW that a particular computation has failed.	{Ref, failure}
Leave Request. Sent by the OW to <code>storage_processes</code> and by <code>storage_processes</code> to <code>store_handlers</code> . When received by a <code>storage_process</code> , it forwards the message to its <code>store_handler</code> . When received by an <code>store_handler</code> , it immediately kills all <code>storage_processes</code> on the node it is running on, and kills itself.	{Pid, Ref, leave}

This table continues on the next page...

Table 1: Description of messages in the system. In this table, we use the initialism “OW” to stand for Outside World: i.e., an arbitrary computer outside of our system.

Message Description	Erlang Pattern
Backup Store Request. Sent by <code>store_handler</code> and <code>storage_process</code> , and received by <code>store_handler</code> . If a <code>store_handler</code> receives this message from a <code>storage_process</code> , it forwards the message to the next <code>store_handler</code> . If an <code>store_handler</code> receives this message from another <code>store_handler</code> , it will back up the data in the message, then notify the OW of the store’s success and the old value.	{Pid, Ref, backup_store, Key, Value, ProcessID}
Joining Behind. A synchronous message received and sent by <code>store_handler</code> . A <code>store_handler</code> will send this when it is joining to the next node’s <code>store_handler</code> to indicate that it is joining behind the recipient in the ring. When received, send all stored backup data to the sender, then delete all backup data for processes numbered less than <code>NodeID</code> . Update local value of <code>prevNodeID</code> , and also change the target node of the Erlang monitoring.	{joining_behind, NodeID}
Joining in Front. Received and sent by <code>store_handler</code> . A <code>store_handler</code> will send this when it is joining to the previous node’s <code>store_handler</code> to indicate that it is joining in front of the recipient in the ring. When receiving such a message, kill the data <code>storage_processes</code> that the new node is now running (i.e. the ones numbered from <code>NodeID</code> to the ID of the node after the new one. Update local value of <code>nextNodeID</code> .	{joining_front, NodeID}
The Erlang node at <code>Node</code> died. Received by <code>store_handler</code> , sent by the Erlang system because the particular <code>store_handler</code> requested it. This is a notification that the node behind the receiving node has stopped running. When such a message is received, the <code>store_handler</code> changes its node’s ID to <code>prevNodeID</code> (the ID of the now-deceased node), which is a data element of the <code>store_handler</code> ’s record. The <code>store_handler</code> then uses all of the backup data it is holding to start up new <code>storage_processes</code> to take over from the node that died. Then it forwards its backup data (of the node that died) to its next node in an <code>appendBackup</code> message, which will now be in charge of backing up that data. Finally it deletes the backup data, and sends a <code>backupRequest</code> message around the ring, to get from the node now behind it the data it should now be backing up, along with that node’s ID.	{nodedown, Node}
Append Backup. Sent and received by <code>store_handlers</code> . When received, the <code>store_handler</code> adds the data contained to its backup data. This is sent when a node <i>A</i> dies, by <i>A</i> ’s successor to <i>A</i> ’s successor’s successor to ensure that the data newly-stored on <i>A</i> ’s successor is backed up. Also updates the value of <code>prevNodeID</code> on <i>A</i> ’s successor’s successor to reflect the new ID of <i>A</i> ’s successor (which is the ID of now-defunct <i>A</i> ).	{appendBackup, BackupData, New-PrevNodeID}

This table continues on the next page...

Table 1: Description of messages in the system. In this table, we use the initialism “OW” to stand for Outside World: i.e., an arbitrary computer outside of our system.

Message Description	Erlang Pattern
Backup Node Data. Received and sent by <code>store_handler</code> . When received, overwrite our backup data with the data in the message; update our <code>prevNodeID</code> , and tell Erlang to monitor the node that’s now previous to us (given in this message as <code>Node</code> ). Sent by a node <i>A</i> ’s predecessor when node <i>A</i> died and <i>A</i> ’s successor is taking over for it, after node <i>A</i> ’s predecessor received a <code>backupRequest</code> message.	{Node, backupNode, NewPrevID, NewBackupData}
Backup Request. Received and sent by <code>store_handler</code> . If it is received on the node whose <code>nextNodeID</code> is <code>DiedNodeID</code> , send each of this node’s <code>storage_processes</code> an <code>all_data</code> message. After compiling all of the results from those requests, send all of this node’s stored data to this node’s successor node in a <code>backupNode</code> message. Otherwise, just forward the request message to the next node’s <code>store_handler</code> . Initially sent by a node which stepped into the void left by a node that died.	{Pid, backupRequest, DiedNodeID}
All data request message. Synchronously received by <code>storage_process</code> and sent by <code>store_handler</code> . When received by a <code>storage_process</code> , respond with a list containing all this <code>storage_process</code> ’s data.	{all_data}
Received by <code>storage_process</code> and sent by <code>store_handler</code> . Simply stops immediately.	terminate