

Distributed Key-Value Storage Algorithm

CS181E — Distributed Systems

Assignment 5

Alejandro Frias

Ravi Kumar

David Scott

April 13, 2014

Algorithm Description

Table 1: Description of messages in the 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 value for key. These are only ever received by storage processes in the network. When an <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 storage process 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 storage process and the data has been backed up in the <code>store_handler</code> sending the message.	{Ref, stored, Old.Value}

This table continues on the next page...

Table 1: Description of messages in the system

Message Description	Erlang Pattern
Retrieve Request. Sent by OW and storage processes; received by storage processes. When a storage processes receives such a message, it checks if the hash of the key is equal to its own process ID. If it is, the storage process has the value for that key, and replies with a retrieve response. If not, it forwards it to the closest storage process to the destination (closest as defined above).	{Pid, Ref, retrieve, Key}
Retrieve Response. Sent by storage processes to the OW. After a storage process receives a retrieve request meant for it, it looks up the relevant value and reports it to the requesting process in the OW.	{Ref, retrieved, Value}
First Key Request. Sent by the OW to storage processes, and by storage processes to storage handlers. When a storage process receives this, it will forward the message up to its storage handler. When such a message is received by a storage handler, 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 next node in the ring.	{Pid, Ref, first_key}
Last Key Request. Sent by the OW to storage processes, and by storage processes to storage handlers. When a storage process receives this, it will forward the message up to its storage handler. When such a message is received by a storage handler, 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 next node in the ring.	{Pid, Ref, last_key}
Num Keys Request. Sent by the OW to storage processes, and by storage processes to storage handlers. When a storage process receives this, it will forward the message up to its storage handler. When such a message is received by a storage handler, 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 next node in the ring.	{Pid, Ref, num_keys}

This table continues on the next page...

Table 1: Description of messages in the system

Message Description	Erlang Pattern
Node List Request. Sent by the OW to storage processes, and by storage processes to storage handlers. 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 query the global registry for the list of nodes and report that data to the requester.	{Pid, Ref, node_list}
Request Result. Sent to the OW by a storage handler. This reports the result of a First Key, Last Key, Num Keys, or Node List Request to the original requester after the storage handlers have finished computing the result.	{Ref, result, Result}
Failure Notification. Sent to the OW by storage handlers or storage processes to notify the OW that a particular computation has failed.	{Ref, failure}
Leave Request. Sent by the OW to storage processes and by storage processes to storage handlers. 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 storage processes on the node it is running on, and kills itself.	{Pid, Ref, leave}
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}

This table continues on the next page...

Table 1: Description of messages in the system

Message Description	Erlang Pattern
<p>Messages About Keys. Sent and received by <code>store_handler</code>. If <code>Ref</code> 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>.</p>	<code>{Pid, Ref, *_key, ComputationSoFar}</code>
<p>Joining Behind. 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>.</p>	<code>{Pid, joining_behind, NodeID}</code>
<p>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 storage processes 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).</p>	<code>{joining_front, NodeID, DestID}</code>
<p>Node with <code>NodeID</code> Died. Received by <code>store_handler</code>, sent by a <code>gen_server</code> listener started by that particular <code>store_handler</code>. 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 the node's ID to <code>NodeID</code>, then uses all of the backup data it's holding to start up new data storage processes. Then it deletes the backup data, and sends a <code>backup_request</code> message around the ring, to get the data it should be backing up from the node behind it.</p>	<code>{died, NodeID}</code>
<p>Backup Node Data. Received and sent by <code>store_handler</code>. When received, add all the data to existing backup data. Send 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.</p>	<code>{backup_node, Data}</code>

This table continues on the next page...

Table 1: Description of messages in the system

Message Description	Erlang Pattern
Backup Request. Received and sent by <code>store_handler</code> . If it is received on the node with <code>DestID</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>backup_node</code> message. If this node is not <code>DestID</code> , 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.	<code>{backup_request, DestID}</code>
All data request message. Received by <code>storage_process</code> and sent by <code>store_handler</code> . When received by a <code>storage_process</code> , respond with an <code>all_data_send</code> message containing all this <code>storage_process</code> 's data.	<code>{all_data, Pid}</code>
All data send message. Received by <code>store_handler</code> and sent by <code>storage_process</code> . The <code>store_handler</code> adds the received data to an ongoing list of data and removes the sender from the list of processes it is waiting for. If it's the last response that was being waited for, send the <code>backup_node</code> message to the next node.	<code>{all_data, Data, Pid}</code>

Our Setup and Some 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 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 does so in the middle of the node that is responsible for the most storage processes (breaking ties arbitrarily by process ID number). 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 needs access to), and in the process tells that next node to delete the first half of its backup data since the new node will 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.

When a node crashes or leaves, the next node finds out, as it 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 we lost some backup data of 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 backup 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. 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 it 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 its 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. These system processes scale with the number of nodes and 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 storage processes deal with completely without the handler. 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.