

Second Assignment Report

SDIS 2020/2021 - MIEIC

Distributed Backup Service for the Internet

Ana Inês Oliveira de Barros

up201806593

Eduardo da Costa Correia

up201806433

João de Jesus Costa

up201806560

João Lucas Silva Martins

up201806436

Introduction

The main goal of this report is to describe the design and implementation of the developed peer-to-peer distributed backup service for the Internet developed for the Distributed Systems course unit.

In our case, we opted to develop on top of what was already done in the first project, so a good foundation was already available for us to focus on the new main features.

Overview

Just like in the first project, we support all the specified operations for the backup service:

- Backup of a file, with a desired replication degree;
- Restore of a backed up file;
- Delete a file's backup;
- Limit local storage used by the application;
- Retrieve service state information.

Summary of features:

- **SSLEngine**: To provide secure communication through TCP connections between peers.
- **Chord**: To guarantee the service's scalability with fault tolerance mechanisms, like:
 - Successor list;
 - Fix fingers;
 - Notify;
 - Stabilize.
- **Thread Pools**: To handle and send multiple messages at once.
- **Java NIO**: To execute non-blocking I/O procedures and handle incoming connections.

Protocols

RMI

RMI provides the mechanism by which the server and the client communicate and pass information back and forth. In our case, both the *Peer* class and the *ChordNode* class act as RMI servers programs which need to create the initial remote objects and export them to the RMI runtime, making them available to receive incoming remote invocations from the client application.

TestApp

Interface implemented by the *Peer* class:

```
4 public interface TestInterface extends Remote {
5     String join() throws RemoteException;
6
7     String backup(String filePath, Integer replicationDegree) throws RemoteException;
8
9     String restore(String filePath) throws RemoteException;
10
11    String delete(String filePath) throws RemoteException;
12
13    String reclaim(int maxCapacity) throws RemoteException;
14
15    String state() throws RemoteException;
```

[TestInterface.java](#)

The following setup procedure creates and exports one remote object and registers it with the RMI registry.

```
552 public static void main(String[] args) {
553     // parse cmdline args
554     Peer prog = null;
555
556     try {
557         prog = new Peer(args);
558     } catch (IOException e) {
559         System.err.println("Couldn't initialize the program.");
560         e.printStackTrace();
561         usage();
562     }
563
564     assert prog != null;
```

[Peer.java](#)

```
576 // setup the access point
577 TestInterface stub;
578 try {
579     stub = (TestInterface) UnicastRemoteObject.exportObject(prog, 0);
580     prog.registry = LocateRegistry.getRegistry();
581     prog.registry.bind(("peer" + prog.id), stub);
582 } catch (Exception e) {
583     System.err.println("Failed setting up the access point for use by the testing app.");
584     System.exit(1);
585     // e.printStackTrace();
586 }
587 try {
588     prog.initCoordNode();
589 } catch (IOException e) {
590     System.err.println("Couldn't create node socket");
591 }
```

[Peer.java](#)

The *Peer* class also provides an implementation for each remote method in the remote interface *TestInterface* shown above. Here is an example of an implementation of one of these methods (*delete*):

```
399 public String delete(String filePath) throws RemoteException {
400     String fileId;
401     try {
402         fileId = DigestFile.getHash(filePath);
403     } catch (IOException e) {
404         throw new RemoteException("Deletion of " + filePath + " failed.");
405     }
406     return "File " + filePath + " deletion: " + this.deleteFromId(fileId);
407 }
```

[Peer.java](#)

For client testing, the *TestApp* class implements the client role by allowing us to invoke the sub protocols provided by the service to back up, restore and delete files, as well as to reclaim the storage space being used by the service and to inspect the internal state of a peer.

In the code snippet below, the client looks up for the remote object, using the same name used by the *Peer* class to bind its remote object. Also, the client uses the *LocateRegistry.getRegistry* API to create a remote reference to the registry on the server's host.

```
142 public void join(ChordInterface nprime) throws RemoteException {
143     this.predecessor = null;
144     this.setSuccessor(nprime.findSuccessor(this.getId()));
145
146     // init finger table
147     for (int i = 0; i < m; ++i) {
148         int fingerStartId = this.getFingerStartId(i);
149         this.fingerTable[i] = nprime.findSuccessor(fingerStartId);
150     }
151 }
```

[chord/ChordNode.java](#)

The client then invokes the lookup method on the registry to look up the remote object by name in the server host's registry. This returns a reference to a registry at the named host and the default registry port (usually 1099). Finally, the client creates a new *TestInterface* object and invokes the desired method of the remote object.

```
19 TestInterface stub = null;
20 try {
21     Registry registry;
22     if (rmiinfoSplit.length > 1)
23         registry = LocateRegistry.getRegistry("localhost", Integer.parseInt(rmiinfoSplit[1]));
24     else
25         registry = LocateRegistry.getRegistry();
26     stub = (TestInterface) registry.lookup(rminame);
27 } catch (RemoteException | NotBoundException e) {
28     System.err.println("Couldn't find/get the desired remote object.");
29     e.printStackTrace();
30     System.exit(1);
31 }
32 assert stub != null;
33
34 String reply = null;
35 String oper = args[1];
36 try {
37     switch (oper.toUpperCase()) {
38         case "BACKUP":
39             if (args.length != 4) usage();
40             String filePath = args[2];
41             int replicationDegree = Integer.parseInt(args[3]);
42             System.out.println("BACKUP " + filePath + " " + replicationDegree);
43
44             reply = stub.backup(filePath, replicationDegree);
45             break;
```

[TestApp.java](#)

Chord

Similarly, RMI provides the mechanism by which a *Peer* joins the Chord ring and by which each *ChordNode* instance communicates with their successors, predecessor, and *fingers*.

Each *ChordNode* instance exposes a few methods to be used by other nodes on the chord ring. Interface that exposes these methods:

```
10 public interface ChordInterface extends Remote {
11     int getId() throws RemoteException;
12     ChordInterface getPredecessor() throws RemoteException;
13     ChordInterface getSuccessor() throws RemoteException;
14
15     ChordInterface[] getSuccessors() throws RemoteException;
16
17     ChordInterface findSuccessor(int id) throws RemoteException;
18     ChordInterface findPredecessor(int id) throws RemoteException;
19     ChordInterface closestPrecedingNode(int id) throws RemoteException;
20
21
22     void notify(ChordInterface n) throws RemoteException;
23
24     InetAddress getAddress() throws RemoteException;
25     int getPort() throws RemoteException;
26     Map<Pair<String, Integer>, Integer> getStoredChunksIds() throws RemoteException;
27 }
```

[chord/ChordInterface.java](#)

The following setup procedure creates and exports one remote object and registers it with the RMI registry.

```
49 ChordInterface stub;
50 try {
51     stub = (ChordInterface) UnicastRemoteObject.exportObject(this, 0);
52     registry.bind(this.id.toString(), stub);
53     System.out.println("Registered node with id: " + this.id);
54     System.out.println(this);
55 } catch (Exception e) {
56     System.err.println("Failed setting up the access point for use by chord node.");
57     e.printStackTrace();
58     System.exit(1);
59 }
```

[chord/ChordNode.java](#)

The *ChordNode* class provides an implementation for each remote method in the remote interface *ChordInterface* shown above. Here is an example of an implementation of one of these methods (*getSuccessor*):

```

82 public ChordInterface getSuccessor() throws RemoteException {
83     boolean goneBad = false;
84
85     for (int i = 0; i < this.succList.length; ++i) {
86         ChordInterface succ = this.succList[i];
87         if (succ == null) break;
88
89         try {
90             succ.getId();
91             if (goneBad) {
92                 this.reconcile(succ);
93
94                 // My successor died, call backup protocol on the chunks i think he was storing
95                 this.backupSuccessorChunks();
96             }
97             return succ;
98         } catch (RemoteException ignored) {
99             goneBad = true;
100             this.succList[i] = null; // node is dead => bye bye
101         }
102     }
103
104     return null;
105
106     // return this.fingerTable[0];
107 }

```

[chord/ChordNode.java](#)

Each node has a join method that automatically joins an existing chord ring, by selecting a random node already there. To do this, the node queries the rmiregistry for the list of objects there and joins the first one available. If no node is found, it is assumed that the node is the only one in the chord ring.

When the lookup is successful, the node initializes its first successor (in the successor list) and its finger table. Both of these are initialized with the assistance of the existing node, by querying who it thinks are successors of the IDs specified (id of the current node and IDs of the fingers).

```

270 ChordInterface node;
271 try {
272     System.out.println(peerId);
273     node = (ChordInterface) this.registry.lookup(peerId);
274 } catch (NotBoundException e) {
275     return "An attempted lookup to a node in the network failed";
276 }
277 this.chordNode.join(node);
278 // Resume pending tasks
279 this.handlePendingTasks();
280 return "Join success";
281 }

```

[Peer.java](#)

Messages and routing

The header of all messages exchanged between Peers follows the following template. In some messages, some fields can be omitted. To signal this, they'll be replaced by “...” in each message.

MESSAGE <file_id> <source_address> <source_port>
<destination_address> <destination_port> <destination_id> <path>*

- **file_id**: ID of the file. Obtained from the name, content, and other metadata of the file using SHA-256 hashing algorithm.
- **source_address**: IPv4 address of the peer which sent the message originally.
- **source_port**: Listening port of peer which sent the message originally.
- **destination_address**: IPv4 address of destination peer.
- **destination_port**: Listening port of destination peer.
- **destination_id**: ID of destination chord node.
- **path**: Collection of chord node IDs through which the message went.

Each Chord node acts as a **router**, so it forwards messages that are not meant for it through the ring to the closest preceding node it is aware of (successor list/finger). For each node the message passes through, the node's ID is added to the path field of the message.

```
89 public void send(Message message) {
90     if (this.messageIsForUs(message)) {
91         System.out.println("\tNot sending message (its for me): " + message + "\n");
92         messageHandler.handleMessage(message);
93     } else { // message isn't for us
94         System.out.println("Sending (ReHopping): " + message + "\n");
95         this.sendToNode(message); // resend it through the chord ring
96     }
97 }
```

[chord/ChordController.java](#)

If a message is meant for the current node (it isn't aware of any preceding node closest to the destination node), then the node will proceed to handle the corresponding message. Otherwise, the current node will send the message to the closest preceding node of the message's destination node (*destination_id* field of the message). In case the current node is the closest one, the message will be sent to the node's successor. This explanation refers to the code below:

```
401 private void sendToNode(Message message) {
402     ChordInterface nextHopDest = null;
403     try {
404         nextHopDest = closestPrecedingNode(message.getDestId());
405     } catch (RemoteException e) {
406         System.err.println("Could not find successor for message " + message + ". Message not sent.");
407         e.printStackTrace();
408     }
409     assert nextHopDest != null;
410     try {
411         if (nextHopDest == this)
412             nextHopDest = getSuccessor();
413         message.setDest(nextHopDest);
414         message.addToPath(nextHopDest.getId());
415     } catch (RemoteException e) {
416         System.err.println("Could connect to chosen next hop dest : TODO Max tries with timeout " + message);
417         e.printStackTrace();
418     }
419     this.sock.send(message);
420 }
421 }
```

[chord/ChordController.java](#)

Backup

The backup service generates an identifier for each file it backs up. It then splits it into chunks and then backs up each chunk independently. To back up a chunk, the initiator-peer sends a PUTCHUNK containing the contents of that chunk.

The PUTCHUNK message is sent to the node responsible for the chunk. The node responsible for the chunk is the node whose ID is the closest to the chunk's ID. PUTCHUNK messages to store replications of this chunk are sent to the next N successors of the responsible node. [PUTCHUNK Message Handler](#).

Header of a PUTCHUNK message:

PUTCHUNK [...] <replication_degree> <sequence_number> <chunk_no>

- **replication_degree**: Desired replication degree of the chunk.
- **sequence_number**: Remaining stores to achieve the desired replication degree.
- **chunk_no**: Number of the chunk being stored.

A peer that stores the chunk upon receiving the PUTCHUNK message, replies with a STORED message to its predecessor. This acknowledgment message is used to update a node's perceived chunks of its successor. The uses and detail implementations of this data structure are specified further in the report. [STORED Message Handler](#)

Header of STORED message:

STORED [...] <chunk_no> <chunk_id>

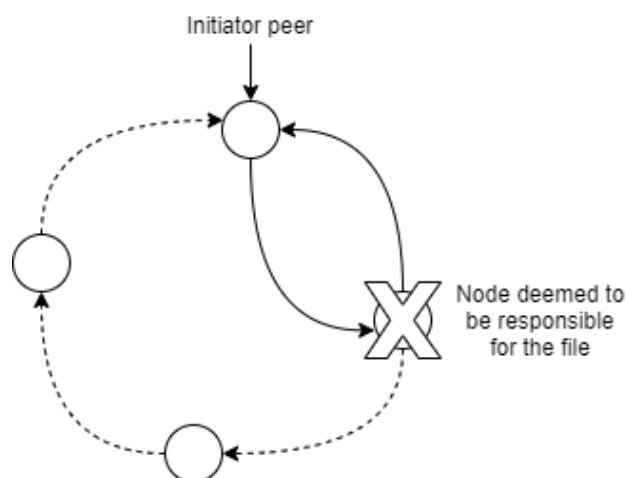
- **chunk_no**: Number of the stored chunk.
- **chunk_id**: ID of the stored chunk.

Restore

The node that starts a restore request of a file will need to recover all the backed up chunks of that file. To recover a chunk, the initiator node sends a GETCHUNK message to the node responsible for that chunk. In case this node hasn't stored the chunk (due to limiting its backup storage space), it will redirect the GETCHUNK message to its successor. The message will be propagated across the chord ring until it reaches a node that is storing the desired chunk. (for example, after all of its space being reclaimed)

In our first implementation of the protocol, the initiator peer was responsible for stopping the propagation of the GETCHUNK. We quickly found that this approach was faulty when handling some cases, such as the following:

1. The initiator node initially sends the message to the node it deems responsible for the file, which corresponds to its direct predecessor.
2. It forwards the message to its successor (back to the initiator), in case it hasn't stored the requested file.
3. The initiator node would wrongly deduce the message had looped through the wrong ring.



Our solution to this problem is to pass the responsibility of stopping the propagation of the message to the responsible peer. This is achieved by using a flag, `has_looped`, which indicates if the message has passed through the responsible node twice. In this case, it will resend the GETCHUNK message directly to the initiator peer. Upon receiving this GETCHUNK message, the initiator peer will be informed that no node in the network is storing the file, by checking if the message was initially sent by him and that flag `has_looped` set to true. If this is the case, it will know that the backup process has failed because at least one chunk wasn't found in the network.

To prevent this, we added two extra fields, *looped* and *responsible*.

[GETCHUNK Message Handler](#)

Header of a GETCHUNK message:

GETCHUNK [...] <chunk_no> <looped> <responsible>

- **chunk_no:** Number of the chunk being retrieved.
- **looped:** True when the message has looped through the chord ring. False, otherwise.
- **responsible:** ID of the node responsible for the chunk

Upon receiving a GETCHUNK message, a node that has a copy of the specified chunk shall send it in the body of a CHUNK message directly to the node who requested the chunk. [CHUNK Message Handler](#)

Header of a CHUNK message:

CHUNK [...] <chunk_no>

- **chunk_no:** Number of the chunk whose body is being sent.

Delete

To delete all file chunks existing in the network, the initiator peer sends a DELETE message to its successor. This message passes through the entire network until it reaches the initiator peer. If a peer has chunks referring to the file specified in the message, they are deleted. [DELETE Message Handler](#)

Header of a DELETE message:

DELETE [...]

Reclaim

When reclaiming disk storage, the node may need to delete a copy of a chunk it has backed up. Upon deleting a backed up chunk, it sends a REMOVED message to the node responsible for the chunk. If the peer that receives this message isn't storing the chunk, the message will be resent through the ring, similarly to the previous messages. When a peer that is either the initiator or is one of the nodes that are storing the chunk file receives the message, it reinitiates the backup protocol. This is achieved by sending a PUTCHUNK with the desired replication degree and respective chunk to the responsible peer of that chunk. [REMOVED Message Handler](#)

Header of a REMOVED message:

REMOVED [...] <chunk_no> <chunk_id>

- **chunk_no**: Number of the removed chunk.
- **chunk_id**: ID of the removed chunk.

State

This operation allows us to observe the service state. In response to such a request, the node sends to the client the following information:

- For each file whose backup it has initiated:
 - The file path name
 - The backup service ID of the file
 - The desired replication degree
 - For each chunk of the file:
 - Its ID

- Its perceived replication degree
- For each chunk it stores:
 - Its ID
 - Its size (in KBytes)
 - The desired replication degree
 - Its perceived replication degree
- For each chunk that the peer's successor is storing:
 - Its fileId
 - Its number
 - Its ID
- The peer's storage capacity, i.e. the maximum amount of disk space that can be used to store chunks, and the amount of storage (both in KBytes) used to back up the chunks.

Example of a peer's state output:

```
CMD: state
Files I initiated the backup of:
Chunks I am storing:
  File ID: d06f4be60f15642ad3b89887fcf829688acb7fe8674c87088482f7f885cb22de
  Chunk no: 1 Id: 7
  Size: 30000
  Desired replication degree: 2
Chunks my succ is storing:
  FileId: d06f4be60f15642ad3b89887fcf829688acb7fe8674c87088482f7f885cb22de ChunkNo: 0 ChunkId: 6
5
  FileId: d06f4be60f15642ad3b89887fcf829688acb7fe8674c87088482f7f885cb22de ChunkNo: 2 ChunkId: 9
6
  FileId: d06f4be60f15642ad3b89887fcf829688acb7fe8674c87088482f7f885cb22de ChunkNo: 1 ChunkId: 7
Storing 30KB of a maximum of infinite KB.
```

Concurrency design

Multi-threading and Thread pools

The *SocketThread* class implements the **Runnable** interface and contains a **fixed thread pool** from the [java.util.concurrent](https://docs.oracle.com/javase/7/docs/api/java/util/concurrent/) package. This type of pool always has a specified number of threads running. Whenever a task terminates, a new task from the thread pool waiting queue is picked up from the newly free thread.

Tasks are submitted to the pool via an internal queue, which holds extra tasks whenever there are more active tasks than threads (waiting queue). Whenever a message is received or needs to be sent, a message handler task is added to the corresponding thread pool, for execution:

```
479 public void send(Message message) {
480     this.sendThreadPool.execute(() -> this.sendInner(message));
481 }
```

[sender/SocketThread.java](#)

```
360 private void readOuter(SelectionKey key, SocketChannel socketChannel, SSLEngineData d) {
361     try {
362         boolean isClosed = this.read(socketChannel, d);
363
364         // create message instance from the received bytes
365         if (isClosed) {
366             key.cancel();
367             d.thread.shutdown();
368
369             if (d.content.size() > 0) {
370                 ByteArrayInputStream bis = new ByteArrayInputStream(d.content.toByteArray());
371                 d.content.reset();
372
373                 ObjectInput in = new ObjectInputStream(bis);
374                 Message msg = (Message) in.readObject();
375                 bis.close();
376
377                 // handle message
378                 this.receiveThreadPool.submit(() -> this.observer.handle(msg));
379             }
380         }
381     } catch (IOException | ClassNotFoundException ignored) {
382         key.cancel();
383         d.thread.shutdown();
384         // System.err.println("Lost message.");
385     }
386 }
```

[sender/SocketThread.java](#)

Appropriate data structures and atomic values

With the use of threads, there is a need for using appropriate data structures and atomic values (in critical regions). We used a [ConcurrentHashMap](#) from the [java.util.concurrent](#) package which supports full concurrency of retrievals and adjustable expected concurrency for updates. We also used [AtomicBooleans](#) from the [java.util.concurrent.atomic](#) package. This package contains a small toolkit of classes that support lock-free thread-safe programming on single variables. The code snippets below illustrate some examples where these classes are used.

```
21 private final AtomicBoolean running = new AtomicBoolean(false);
22 private final ExecutorService sendThreadPool = Executors.newFixedThreadPool(MAX_CONNS);
```

[src/sender/SockThread.java](#)

```
19 InetAddress address;
20 Integer port;
21 private final ConcurrentMap<Pair<String, Integer>, CompletableFuture<byte[]>> receivedChunks;
22
23 public MessageHandler(SockThread sock, ChordController chordController) {
24     this.controller = chordController;
25     this.receivedChunks = new ConcurrentHashMap<>();
26     this.address = sock.getAddress();
27     this.port = sock.getPort();
28 }
```

[src/sender/MessageHandler.java](#)

Synchronized methods and statements

Synchronized methods and statements enable a simple strategy for preventing thread interference and memory consistency errors. We used these in our project to achieve better concurrency:

```
136 synchronized (State.st) {
137     DigestFile.deleteFile(message.getFileId());
138     State.st.removeSuccChunk(message.getFileId());
139 }
```

[sender/MessageHandler.java](#)

Java NIO

The Java NIO Selector is a component which can examine one or more Java NIO Channel instances, and determine which channels are ready for operations like reading and/or writing. This way, a single thread can monitor multiple channels, and thus multiple network connections.

In our project, Java NIO is used for the writing and reading of messages, and accepting of connections. Each peer acts both as a server and as a client. A selector is used to accept new connections and read from existing sockets. All socket channel reading and writing operations are asynchronous.

```
334 Iterator<SelectionKey> selectedKeys = this.selector.selectedKeys().iterator();
335 while (selectedKeys.hasNext()) {
336     SelectionKey key = selectedKeys.next();
337     selectedKeys.remove();
338     if (!key.isValid())
339         continue;
340
341     try {
342         if (key.isAcceptable()) {
343             this.acceptCon(key);
344         } else if (key.isReadable()) {
345             SocketChannel socketChannel = (SocketChannel) key.channel();
346             SSLEngineData d = (SSLEngineData) key.attachment();
347
348             try {
349                 d.thread.submit(() -> this.readOuter(key, socketChannel, d));
350             } catch (RejectedExecutionException ignored) {
351                 key.cancel();
352             }
353         }
354     } catch (CancelledKeyException ignored) {
355     }
356 }
```

[sender/SockThread.java](#)

Upon accepting a new connection, we attribute the new socket channel a new thread (up to a limit). Every time this socket channel is notified about pending read operations, it uses this thread (single thread thread-pool). This is important, because it allows us to read concurrently and guarantees that the messages of each SSLEngine are processed sequentially.

JSSE

Peers send and receive messages through TCP socket channels. The data flow through these channels is encrypted using Java's `SSLEngine`. This means that `SSLEngine` is used for all message communication (all protocols).

Our implementation requires client authentication, JKS, and SunX509. TLS version 1.3 is used, and our code conforms to the half-close policy (new requirement after TLS version 1.2). This means that all connections end with the independent closure of both the outbound and inbound ends, thus preventing truncation attacks.

Both reading and writing is done with non-blocking sockets and Java NIO's features (discussed further in the scalability section of the report). This is the main reason why `SSLEngine` is used.

```
195 hs = res.getHandshakeStatus();
196
197 // check status
198 switch (res.getStatus()) {
199     case OK:
200         // do nothing ?
201         break;
202     case BUFFER_OVERFLOW:
203         // the client maximum fragment size config does not work?
204         peerAppData = this.handleOverflow(engine, peerAppData);
205         break;
206     case BUFFER_UNDERFLOW:
207         // bad packet, or the client maximum fragment size config does not work?
208         // we can increase the size
209         peerNetData = this.handleUnderflow(engine, peerNetData, peerAppData);
210         // Obtain more inbound network data for src,
211         // then retry the operation.
212         break;
213     case CLOSED:
214         engine.closeOutbound();
215         return 0;
216     default:
217         throw new IllegalStateException("Unexpected value: " + res.getStatus());
218 }
```

[sender/SockThread.java](#)

After the server detects that the connection is being closed by the peer (through the unwrapping of a `CLOSE` message), it proceeds with an orderly shutdown. If the connection is shut down by the peer without following the orderly shutdown procedure, the server still attempts to close the connection in an orderly fashion. The received message is processed after the connection is closed, even in the case of *end of stream* errors or other non-orderly shutdowns.

The client can time out while waiting for the server to acknowledge the orderly shutdown. In these cases, the client closes its socket channel and leaves.

Note: We are using the `SSLEngine` configuration that allows for bigger packets.

Scalability

In order to achieve higher scalability in our project, we decided to implement scalability both at design level and at implementation level. At design level, we implement a **Chord** architecture that provides routing of the messages. At implementation level, we used thread-pools to simultaneously handle several processes and Java NIO to achieve asynchronous I/O operations. All of these subjects have already been described in previous sections in this report.

Furthermore, to prevent hashing collision of both file and Chord node's IDs, we used the SHA-256 algorithm.

Additionally, a node may join the ring directly, by specifying the IP address and port of an existing node in the ring. Otherwise, it is also possible for a node to automatically join any of the nodes in the chord ring by looking at the objects that are bound to the RMI registry. When joining the chord ring, the new node queries the existing node for its successor and fingers, thus joining the ring.

Fault-tolerance

Since we implemented the Chord protocol, we took advantage of its fault-tolerant features to achieve a robust and reliable implementation.

We are using $m = 128$, so we have 7 fingers in our finger table. Assuming each node has a 50% chance to fail (quite high!), we achieve the following chances of all nodes failing:

- 3 successors: $\approx 12\%$
- 4 successors: $\approx 6\%$
- 5 successors: $\approx 3\%$
- 6 successors: $\approx 1\%$
- **7 successors: $\approx 0\%$**

With this in mind, we decided to go with the **7 successors list size**, since it is virtually impossible for all the successors of a node to fail at once.

We use the **SHA-256 hashing function** to generate IDs. This helps us achieve some tolerance to denial-of-service (DoS) attacks.

We call the *fixFingers* method periodically to refresh the finger table entries:

```
213 public void fixFingers() {
214     if (nextFingerToFix >= m)
215         nextFingerToFix = 0;
216
217     int succId = this.getFingerStartId(nextFingerToFix);
218     try {
219         // System.out.println("I am " + this.id + " and I'm updating finger " + nextFingerToFix + ", id: " +
succId);
220         fingerTable[nextFingerToFix] = this.findSuccessor(succId);
221         // System.out.println("They tell me it's: " + fingerTable[nextFingerToFix].getId());
222     } catch (RemoteException e) {
223         fingerTable[nextFingerToFix] = this;
224     }
225     ++nextFingerToFix;
226 }
```

[chord/ChordNode.java](#)

We use the *notify* method to notify a node's successor of its existence, so the successor can change its predecessor to this node.

```
193 public void notify(ChordInterface nprime) throws RemoteException {
194     int nprimeId = nprime.getId();
195
196     if (this.predecessor == null) {
197         this.predecessor = nprime;
198     } else {
199         try {
200             int predecessorId = this.predecessor.getId();
201             if (ChordNode.inBetween(nprimeId, predecessorId, this.id))
202                 this.predecessor = nprime;
203         } catch (RemoteException ignored) {
204             // if our predecessor died, we accept the new one
205             this.predecessor = nprime;
206         }
207     }
208 }
```

[chord/ChordNode.java](#)

To stabilize the network and guarantee that it is consistent, the *stabilize* method is called periodically. The corresponding node asks the successor about its predecessor, verifies if its immediate successor is accordingly and tells the successor about it.

```
164 public void stabilize() throws RemoteException {
165     // TODO if statement is needed?
166     //if (this.predecessor != null) {
167     //}
168
169     ChordInterface succ = this.getSuccessor();
170     if (succ == null) return; // very bad
171
172     // update predecessor
173     ChordInterface me$ = succ.getPredecessor();
174     if (me$ != null) { // if pred knows about a succ
175         try {
176             if (ChordNode.inBetween(me$.getId(), this.id, succ.getId())) {
177                 this.setSuccessor(me$);
178                 succ = me$;
179             }
180         } catch (RemoteException ignored) {
181             // if we fail getting the me$ id, it is dead and we don't want it as a successor
182         }
183     }
184     succ.notify(this); // notify successor about us
185
186     this.reconcile(succ);
187 }
```

[chord/ChordNode.java](#)

The *checkPredecessor* method periodically checks if a predecessor is alive.

```
322 public void checkPredecessor() {
323     if (predecessor != null) {
324         try {
325             this.predecessor.getId();
326         } catch (RemoteException e) {
327             this.predecessor = null;
328         }
329     }
330 }
```

[chord/ChordNode.java](#)

When a peer disconnects the network (due to crashing or simply quitting), its predecessor will send a REMOVED message as described in the RECLAIM sub protocol.

This will ensure that the replication degree will be maintained in the network, as the backup protocol will be reinitiated again if a node that is either storing the file or is its initiator exists. To achieve this, when a node sets its new direct successor, it stores a record of all the

new successors' chunk IDs. Furthermore, when a node stores a new file or deletes an old one, it sends its predecessor either a STORED or a REMOVED message, respectively.

This approach greatly improves fault tolerance since the only way that a stored chunk can be lost in the network is if each pair of predecessors and successors crash simultaneously, which is highly unlikely.

Conclusion

To conclude, we think that this project allowed us to better understand the adequate protocols and algorithms used to develop secure, fault tolerant and scalable systems on the internet, by applying the knowledge acquired in the courses classes. It also allowed us to refine some concepts applied in the first project and implement our own solution.

That being said, we designed this project to be as robust as we could make it and handle all kinds of possible scenarios, leading us to believe we implemented it successfully and achieved the goals of this course unit.

Despite this, we believe that the protocol could have been implemented in a slightly different way to achieve better results and performance.

Firstly, we could use a network overlay like Pastry instead of Chord, because it works better on file/chunk backup systems that support replication.

Secondly, we could send only the headers of the messages and make it so the peers that want the bodies of them, ask for it directly from the initiator peer. Although doing this would make the protocol more complex, it would reduce the size of the messages that are sent in each hop, reducing the stress on the network.

Thirdly, we could make use of an input buffer to save incoming data. This would allow us to send arbitrarily sized messages and, consequently, divide the files into bigger chunks (currently the files are divided into chunks of 30 KB). These changes would require some minor design changes relating to the way the data is read, but the benefits would be noticeable. Also, we would need to divide the messages that we're sending.

Bibliography

1. Ion Stoica, Robert Morris, David Karger, M. Frans Kaashoek, Hari Balakrishnan. 2001. *Chord: A Scalable Peer-to-peer Lookup Service for Internet*.
2. Oracle. "Java Secure Socket Extension (JSSE) Reference Guide".
<https://docs.oracle.com/en/java/javase/15/security/java-secure-socket-extension-jsse-reference-guide.html>
3. alkarn. "Server and Client implementation with SSLEngine".
<https://github.com/alkarn/sslengine.example>
4. Maarten van Steen, Andrew S. Tanenbaum. 2018. *Distributed Systems*.
5. Olliotte Rusty Harold. 2013. *Java Network Programming*.