

# ASCENS

## Autonomic Service-Component Ensembles

### TR-1202: Science Cloud: Modelling and Implementing the Peer-to-Peer DHT protocol 'Chord'

Grant agreement number: **257414**  
Funding Scheme: **FET Proactive**  
Project Type: **Integrated Project**  
Latest version of Annex I: **7.6.2010**

Author(s): **Philipp Zormeier (LMU), Annabelle Klarl (LMU), Christian Kroiß (LMU), Philip Mayer (LMU)**

Periodic report:  
Periodic covered:  
Date of technical report: **July 25, 2012**  
Revision: **V1**  
Classification: **[CO]**

Project coordinator: **Martin Wirsing (LMU)**  
Tel: **+49 89 2180 9154**  
Fax: **+49 89 2180 9175**  
E-mail: **wirsing@lmu.de**

Partners: **LMU, UNIPI, UDF, Fraunhofer, UJF-Verimag, UNIMORE, ULB, EPFL, VW, Zimory, UL, IMT, Mobsya, CUNI**



## Abstract

The ASCENS project is an EU-funded integrated project within the 7th framework program with the goal of handling open-ended, highly parallel, massively distributed systems by breaking them down into self-aware, self-adaptive and self-expressive autonomic components which are grouped into so-called *ensembles*.

One of the case studies to evaluate the ASCENS results is the *Science Cloud* case study which describes a peer-to-peer voluntary computing cloud for scientific usage which consist of an arbitrary number of nodes running the science cloud software. The nodes interact to ensure that applications deployed in the cloud are kept running in case of disappearing nodes, changing load conditions, and different cloud layouts.

One possible option for implementing communication within a distributed peer-to-peer system is the Chord protocol which induces an overlay ring-like structure within the cloud to ensure both resilience and fast access to values in the node-based distributed hash table.

In this document, we describe both the modelling and the implementation of Chord as a test-case for the usability within the Science Cloud case study. Modelling of Chord is achieved using SCEL (Service Component Ensemble Language), the primary formal modelling language of ASCENS. The implementation of Chord is created using the Java programming language and is kept close to the SCEL model.

## Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>About Chord</b>	<b>5</b>
2.1	Excluding and including predecessors . . . . .	6
2.2	Leaving nodes . . . . .	7
2.3	Data management . . . . .	7
<b>3</b>	<b>About SCEL</b>	<b>8</b>
3.1	Used SCEL dialect . . . . .	8
3.2	Node properties . . . . .	9
3.3	SCEL macros . . . . .	10
<b>4</b>	<b>Modelling Chord in SCEL</b>	<b>11</b>
4.1	Node lookup procedures . . . . .	11
4.2	Node Joins . . . . .	15
4.3	Leaving nodes . . . . .	18
4.4	Data Lookup . . . . .	20
4.5	Data Storage . . . . .	21
4.6	Stabilization . . . . .	21
4.7	Main Node Process . . . . .	23
<b>5</b>	<b>Implementing Chord</b>	<b>23</b>
5.1	Architecture . . . . .	23
5.1.1	Core layer . . . . .	24
5.1.2	Communication layer . . . . .	24
5.1.3	Graphical User Interface . . . . .	25
5.2	Monitoring the System . . . . .	26
<b>6</b>	<b>Conclusion</b>	<b>27</b>



## 1 Introduction

The ASCENS *Science Cloud* case study [vRA<sup>+</sup>11] realizes a peer-to-peer platform-as-a-service by an arrangement of autonomic, self-aware service components. A service component, running on a physical or virtual machine, continually monitors its own status as well as the status of other components associated to the cloud. Components may also be grouped into service component ensembles based on certain properties of the individual components or the links between them. The main aim of the science cloud is ensuring that deployed applications are executed reliably and while keeping to their SLA, even if individual service components fail or disappear.

One option for realizing distributed communication between the individual components of the Science Cloud is the *Chord* protocol [SMLN<sup>+</sup>03]. Chord is a distributed lookup protocol providing simple ID assignment to both nodes and data and ensures efficient lookups. We have extended the basic protocol for our purposes with regard to predecessor handling, leaving nodes, and data handling.

In this document, we formalize the basic Chord protocol with SCEL [GLPT12], a language for modelling autonomic service components and their ensembles, and provide an implementation of Chord in Java as a proof of its practical use.

This document is structured as follows. We first give a short introduction to the Chord protocol in section 2, followed by a description of SCEL in section 3. Afterwards, we discuss our main findings, i.e. the modelling of Chord in SCEL (section 4) and the Java implementation (section 5). We conclude in section 6.

## 2 About Chord

The Chord protocol implements a distributed lookup protocol for locating nodes which store a certain piece of data. Chord thus provides a mapping of some data key to a node which is responsible for this data. The protocol therefore implements a distributed hash table.

The members of a Chord network are logically arranged in a circle which is described by  $m$ -bit addresses. A node is identified by an address  $n$  between 0 and  $2^m - 1$ . The address is generated by hashing its IP address.

Data is then stored at the node to which a key associated to the data item maps (key/value pair). If a data item with key  $m$  is queried along the network, the next node  $n$  in clockwise direction called *successor*( $m$ ) (possibly  $m$ ) is responsible for locating the data.

To establish communication around the network, the independent nodes have to be aware about a couple of other nodes. In a so called *finger table*, each node  $n$  stores information about  $m$  other subsequent successors. *finger*[ $i$ ].*node* is *successor*( $n + 2^{i-1}$ ) ( $i \in \{1, \dots, m\}$ ). *finger*[1].*node* therefore is the direct successor of  $n$ . Though for correct behavior only information about the next successor is needed, the finger table ensures efficient communication (logarithmic complexity). In addition, the main procedures assume knowledge of the node's predecessor.

If any node needs to access data with key  $m$ , it looks up the nearest element in its finger table. This node either is directly responsible for the key or passes it to its own nearest finger.

A joining node initializes its finger table by means of an existing arbitrary node in the network (which must be known). That node locates all  $m$  successors via its own fingers.

After the finger table is set, the new node updates all preceding nodes which might contain it in their finger tables. This procedure ensures valid lookup functionality.

Note that leaving nodes are not specified by the basic protocol. To react to node failures, concurrent events and other challenges, further techniques have to be applied. A Chord-based application might, for example, establish data redundancy and mechanisms that maintain lookup information.

For our model and implementation of Chord, we had to add some enhancements to the original Chord protocol. First, we will describe an inconsistency of the protocol observed during implementation (subsection 2.1). Afterwards we discuss new functionality: Leaving nodes in subsection 2.2 and data handling in subsection 2.3.

## 2.1 Excluding and including predecessors

While testing the implementation we noticed an inconsistency in the Chord protocol, which may lead to invalid finger tables. This situation may occur if a node joins or leaves the network whose id matches the start id of a finger of a preceding node. We will discuss the joining case here as our leaving behavior is yet to be defined.

The procedure *update\_others* notifies nodes which might have to update their finger table with the new node with a call of *update\_finger\_table*.

```
n.update_others()
  for (i = 1 to m)
    p = find_predecessor(n - 2i-1);
    p.update_finger_table(n,i);
```

To determine the node that should check its i-th finger table entry, the predecessor of id  $n - 2^i$  is calculated. This node then eventually updates the entry and passes the message to its direct predecessor.

```
n.update_finger_table(s,i)
  if (s ∈ [n,finger[i].node))
    finger[i].node = s;
    p = predecessor;
    p.update_finger_table(s,i);
```

Within *find\_predecessor*, a predecessor of an id is determined excluding the id itself. If the predecessor of id  $n$  is required and there actually is a node with id  $n$ , the predecessor of that node will be returned.

```
n.find_predecessor(id)
  n' = n;
  while (id ∉ (n', n'.successor])
    n' = n'.closest_preceding_finger(id)
  return n';
```

If  $n'.successor$  is matching the given id, the while loop stops and  $n'$  instead of  $n'.successor$  is returned. If this is the case, however, the procedure *update\_others* will skip the node whose finger  $i$  is already (perfectly) pointing to the new node.

To solve this issue, we provide a variation of *find\_predecessor*. This variation simply checks if  $n'.successor$  is a perfect match.

```
n.find_predecessor_incl(id)
  n' = n;
  while (id ∉ (n', n'.successor])
    n' = n'.closest_preceding_finger(id)
  if (id == n'.successor)
```

```

    return n'.successor;
return n';

```

This new procedure is called instead of *find\_predecessor* in the join and leave procedures. In all other cases the original procedure is adequate.

The other node lookup and joining behavior was not modified but enhanced by data management, which is described below.

We use the *stabilize* and *fix\_finger* procedure as defined in the Chord paper, the only difference being that we do not limit the join procedure to setting the successor as proposed there.

## 2.2 Leaving nodes

We define the following leaving behavior:

```

n.leave()
  finger[1].node.add_data(dataMap);
  finger[1].node.predecessor = predecessor;
  for (i = 1 to m)
    p = find_predecessor(n - 2i-1);
    p.clean_finger(n,i,finger[1].node);

n.clean_finger(s,i,c)
  if (finger[i].node == s)
    finger[i].node = c;
    p = predecessor;
    p.clean_finger(s,i,c);

```

A leaving node updates the predecessor of its successor and afterwards notifies preceding nodes which might have to update their finger table. This behavior is analog to the behavior in *update\_others* in *join*, the only difference being a candidate node being passed along with the message. The candidate to replace the original finger node is its successor.

## 2.3 Data management

As mentioned above, we use Chord to implement a map of key ids to data. This can be done by simply adding local hash maps to the nodes. To transfer keys between the nodes during joins and leaves, these maps have to offer methods like *put*, *get*, *remove*, *extract* and *merge*.

If data at a given key is requested, the successor node of the key is responsible for the data. This leads to the following procedures:

```

n.lookup_data(key)
  if (key ∈ dataMap.domain())
    return dataMap.get(key)
  else
    n' = find_successor(key);
    return n'.lookup_data(key);

n.store_data(key,data)
  if (key ∈ dataMap.domain())
    dataMap.put(key,data)
  else

```

```
n' = find_successor(key);
n'.store_data(key,data);
```

These procedures check if the given id is within the local key space, which is similar to the one of the data map. In this case, the data is returned or stored. In the other case, the responsible node is calculated and the procedure is started there.

To maintain the right data assignment during ongoing joins and leaves, we require the following modifications to *join* and *leave*. In *join*, after the call of *update\_others*, *finger[1].node.extract\_data(predecessor,n)* is called and the resulting map is the new local data map:

```
n.extract_data(p,n')
return dataMap.extract((p,n'));
```

The procedure *extract* is meant to return a map containing all key value pairs within the given interval and shortens the domain by the interval.

When a node leaves, it first has to transmit its data to its successor. For this reason, the whole dataMap is sent to the successor node:

```
n.transmit_data(otherMap)
dataMap.merge(otherMap);
```

The procedure *merge* extends the data map's domain and adds the key-value pairs of the given map. At the beginning of *leave*, *finger[1].node.transmit\_data(dataMap)* has to be called.

All procedures in full length will be repeated at the corresponding place in section 4.

### 3 About SCEL

SCEL (Service Component Ensemble Language) [GLPT12] is a language for modelling autonomic service components and service component ensembles. SCEL is a formal language and built on top of solid semantic grounds. SCEL is also a parameterized language, i.e. the first step when using SCEL is settling on a concrete SCEL *dialect* by substantiating the variable parts of the language (cf. [dNFLP11]):

- *The language for policies* together with an interaction and authorization predicate.
- *The language for representing knowledge items and repositories*, i.e. handling data within the language.
- *The expression language* which is used for producing values and evaluating data.

In the following subsections, we discuss the SCEL dialect used (section 3.1), which data (or rather, knowledge) we store within the nodes (section 3.2), and a set of macros we use to ease the SCEL specification (section 3.3).

#### 3.1 Used SCEL dialect

In order to model in SCEL we define a dialect, i.e. concrete instantiations for the three languages discussed above. We set the stage as follows:

- Our model of the Chord algorithm does not require *policies* since all logics and constraints are directly included within the modelled processes. Consequently, no policy language is required.



- The *knowledge* repositories of nodes consist of tuples. We use a key-based approach to retrieving tuples, i.e. the first element of each tuple contains a key token which is used like an identifier. The rest of the tuple is an arbitrary number of arguments, which vary based on the key or the purpose of the tuple. Tuples are added to the repository by a *put* statement, which does not overwrite tuples with the same key. A *qry* statement returns the first found tuple, and a *get* statement additionally removes the retrieved tuple of the repository. Querying and getting tuples is realized by pattern matching. The statements have to consist of a specific key and the right number of arguments. The other arguments may be certain values or define a pattern to match tuples. Alternatively, a process can bind data to local variables by inserting wild card-like expressions, for example *!x*. These expressions can be extended by patterns in brackets, for example *!x[x > 20]*. Only certain values and patterns narrow the space of results.
- *Expressions* are mainly used within tuples and can be of arbitrary data types. The tuple keys are strings, but other values might also be boolean, numeric or of other types. We use Java-like syntax within curled brackets to symbolize object manipulation. Algebraic expressions are written in plain mathematical syntax and are wrapped by curled brackets.

All other properties and syntax elements are standard SCEL as explained in [GLPT12].

### 3.2 Node properties

Every node knows its finger table, the address of its predecessor and the data it is responsible for. This information is stored in their knowledge repository as follows:

- *dataMap* - a *Map<int,data>* object containing all node data.  
methods:
  - *put(key,data)* - adds the specified key-value pair.
  - *get(key)* - returns the requested value.
  - *domain()* - returns the interval of all assigned keys.
  - *merge(Map<int,data>)* - merges the domains and adds all data of the given map.
  - *extract(domain)* - shortens the domain by the given domain and returns a map containing lost key-values.
  - *addAll(Map<int,data>)* - adds all elements.
- *fingers* - the finger list.  
fields of the fingers:
  - *start* - the expected node position.  
The *i*-th finger has  $start = n + 2^{i-1} \bmod 2^m$ , *n* being the nodes identifier.
  - *node* - *successor(start)*.
 methods:
  - *set(int,finger)* - sets given finger.
  - *get(int)* - returns an element.
- *pred* - the node's predecessor

### 3.3 SCEL macros

In order to abbreviate the specification of the SCEL model we define the following three macros.

- *Sequential composition*

To compose two processes sequentially, we use the following syntax:

$$P \triangleq P_1 \cdot P_2$$

Sequential composition can be realised by synchronization over the knowledge repository. Here is an example:

$$P_1 \cdot P_2 \triangleq P_1^*[P_2^*] \text{ with an empty policy and}$$

$$\begin{aligned} P_1^* &\triangleq \mathbf{put}(\text{"uniqueIdentifier1"}, \{\text{generateID}()\})@self \\ &\quad .\mathbf{get}(\text{"uniqueIdentifier1"}, !id)@self \\ &\quad .\mathbf{put}(\text{"uniqueIdentifier2"}, id)@self \\ &\quad .P_1 \\ &\quad .\mathbf{put}(\text{"ready"}, id)@self \\ &\quad .nil \\ P_2^* &\triangleq \mathbf{get}(\text{"uniqueIdentifier2"}, !id)@self \\ &\quad .\mathbf{get}(\text{"ready"}, id)@self \\ &\quad .P_2 \end{aligned}$$

Process  $P_1$  has to end with nil. By writing  $P_1.\mathbf{put}(\text{"ready"}, id)@self.nil$ , we mean to add the put action right before  $P_1$ 's nil.

- *Loops*

$(i : start \rightarrow end)(P(i), Q)$  abbreviates

$$\begin{aligned} &\dots.\mathbf{put}(\text{"i"}, start)@self \\ &\quad .\mathbf{For} \\ \mathbf{For} &\triangleq (\mathbf{Check} \cdot P(i) \cdot \mathbf{Increment}) + \mathbf{Finally} \\ \mathbf{Check} &\triangleq \mathbf{get}(\text{"i"}, !i[|start - i| \leq |start - end|]) \\ &\quad .nil \\ \mathbf{Increment} &\triangleq \mathbf{put}(\text{"i"}, \{i + (end - start)/|end - start|\})@self \\ &\quad .\mathbf{For} \\ \mathbf{Finally} &\triangleq \mathbf{get}(\text{"i"}, !i[|start - i| > |start - end|])@self \\ &\quad .Q \end{aligned}$$

$P(i)$  has to with process nil and  $Q$  is a process. This macro simplifies the expression of for loops. An equivalent Java example would be

```
for(int i = start; i < end, i ++){
    P(i)
}
Q
```

As usual, we can use standard comparison operators for checking the loop end condition, i.e.  $<$ ,  $\leq$ ,  $>$ , and  $\geq$ , and both increment  $++$  and decrement  $--$  the loop variable.

- *Parallel Processes*

For readability we write  $P_1|P_2|\dots|P_n$ , meaning  $P_1[P_2[\dots[P_n]\dots]]$  with empty policies.

## 4 Modelling Chord in SCEL

This section gives a model of Chord using the SCEL language as parameterized in the previous section. In the following, we assume that value  $m$  is known by every node in the network. Please note that summation is considered modulo  $2^m$  and *self* represents both the processing node and its identifier.

### 4.1 Node lookup procedures

The Chord protocol contains the following lookup procedures:

```

n.find_successor(id)
  n' = find_predecessor_excl(id);
  return n'.successor;

n.find_predecessor_excl(id)
  n' = n;
  while (id ∉ (n', n'.successor])
    n' = n'.closest_preceding_finger(id)
  return n';

n.find_predecessor_incl(id)
  n' = n;
  while (id ∉ (n', n'.successor])
    n' = n'.closest_preceding_finger(id)
  if (id == n'.successor)
    return n'.successor;
  return n';

n.closest_preceding_finger(id)
  for (i = m downto 1)
    if (finger[i].node ∈ (n,id))
      return finger[i].node;
  return n;

```

We now discuss the SCEL model for each of these.

#### **find\_successor**

Process findSucc waits for successor requests and calculates the response according to the Chord protocol. As this is the first SCEL process shown, we describe our modelling approach here. Within this model, procedures, which can be called locally or remotely, are processes that wait for calls.

Those calls take place via the local knowledge repositories, as well the result return mechanisms. So this `findSucc` waits for calls by the blocking `get..` statement. A calling process has to *put* a tuple with the key “`findSuccReq`”, its id for returns and the target id as arguments in the knowledge repository of the system running this `findSucc` process, to trigger calculation. The given parameters are then stored within the local variables `reqNode` and `targID`. In the following procedure `find_predecessor` is called by triggering process `findPred` with an analog pattern. The result is retrieved via the knowledge repository. The last but one line shows that returning values is indeed handled by putting them in knowledge repositories. The following process is `findSucc` itself so the procedure `find_successor` remains available.

```
findSucc  $\triangleq$  get(“findSuccReq”, !reqNode, !targID)@self
      .put(“findPredExclReq”, self, targID)@self
      .get(“findPredExclResp”, targID, !targPredNode)@self
      .put(“getSuccReq”, self)@targPredNode
      .get(“getSuccResp”, targPredNode, !resultNode)@self
      .put(“findSuccResp”, targID, resultNode)@reqNode
      .findSucc
```

### **get\_successor**

We add the process `getSucc` to access other nodes' successors. It realizes a getter which encapsulates finger table access.

```
getSucc  $\triangleq$  get(“getSuccReq”, !reqNode)@self
      .qry(“fingerTable”, !fingers)@self
      .put(“getSuccResp”, self, {fingers.get(1).node})@reqNode
      .getSucc
```

### **find\_predecessor\_excl**

Procedure `find_predecessor_excl` contains a while statement, which can be expressed via nondeterministic choices with guards.

```
findPredExcl  $\triangleq$  get(“findPredExclReq”, !reqNode, !targID)@self
      .put(“tmpPredNode”, self)@self
      .findPredExclWhile(reqNode, targID)
```

```
findPredExclWhile(reqNode, targID)  $\triangleq$  get(“tmpPredNode”, !tmpPredNode)@self
      .put(“getSuccReq”, self)@tmpPredNode
      .(findPredExclWhileBody(reqNode, targID, tmpPredNode)
      +
      findPredExclResp(reqNode, targID, tmpPredNode))
```

The process, that wins the nondeterministic choice, is the one that starts with the matching *get* statement. The other process at the same time dies.

```

findPredExclWhileBody(reqNode,
  targID, tmpPredNode)  $\triangleq$  get("getSuccResp", tmpPredNode, !tmpPredSuccNode
  [targID  $\notin$  (tmpPredNode, tmpPredSuccNode)])@self
  .put("predFingerReq", self, targID)@tmpPredNode
  .get("predFingerResp", targID, !tmpPredNode)@self
  .put("tmpPredNode", tmpPredNode)@self
  .findPredExclWhile(reqNode, targID)

```

```

findPredExclResp(reqNode,
  targID, tmpPredNode)  $\triangleq$  get("getSuccResp", tmpPredNode, !tmpPredSuccNode
  [targID  $\in$  (tmpPredNode, tmpPredSuccNode)])@self
  .put("findPredExclResp", targID, tmpPredNode)@reqNode
  .findPredExcl

```

### **find\_predecessor\_incl**

This procedure is defined analogue to *find\_predecessor\_excl*, but with our modification described in 2.1.

```

findPredIncl  $\triangleq$  get("findPredInclReq", !reqNode, !targID)@self
  .put("tmpPredNode", self)@self
  .findPredInclWhile(reqNode, targID)

```

```

findPredInclWhile(reqNode, targID)  $\triangleq$  get("tmpPredNode", !tmpPredNode)@self
  .put("getSuccReq", self)@tmpPredNode
  .(findPredInclWhileBody(reqNode, targID, tmpPredNode)
  +
  findPredInclResp(reqNode, targID, tmpPredNode))

```

```

findPredInclWhileBody(reqNode,
  targID, tmpPredNode)  $\triangleq$  get("getSuccResp", tmpPredNode, !tmpPredSuccNode
  [targID  $\notin$  (tmpPredNode, tmpPredSuccNode)])@self
  .put("predFingerReq", self, targID)@tmpPredNode
  .get("predFingerResp", targID, !tmpPredNode)@self
  .put("tmpPredNode", tmpPredNode)@self
  .findPredInclWhile(reqNode, targID)

```

```

findPredInclResp(reqNode,
  targID, tmpPredNode)  $\triangleq$  get("getSuccResp", tmpPredNode, !tmpPredSuccNode
    [targID  $\in$  (tmpPredNode, tmpPredSuccNode)])@self
  .put("tmpPredSuccNode", tmpPredSuccNode)@self
  .(findPredInclResult(reqNode, targID, tmpPredNode)
    +
    findPredInclRespResultSucc(reqNode, targID, tmpPredNode))

```

```

findPredInclRespResult(reqNode,
  targID, tmpPredNode)  $\triangleq$  get("tmpPredSuccNode", !tmpPredSuccNode
    [targID! = tmpPredSuccNode])@self
  .put("findPredInclResp", targID, tmpPredNode)@reqNode
  .findPredIncl

```

```

findPredInclResultSucc(reqNode,
  targID, tmpPredNode)  $\triangleq$  get("tmpPredSuccNode", !tmpPredSuccNode
    [targID == tmpPredSuccNode])@self
  .put("findPredInclResp", targID, tmpPredSuccNode)@reqNode
  .findPredIncl

```

### closest\_preceding\_finger

predFinger implements *closest\_preceding\_finger*. This is the first process that contains the loop macro.

```

predFinger  $\triangleq$  get("predFingerReq", !reqNode, !targID)@self
  .put("break", false)@self
  .(i : m  $\rightarrow$  1)(predFingerFor(i, reqNode, targID),
    predFingerReturn(reqNode, targID))

```

Each iteration begins with a check of the *break* variable. If *break* is set to true, the result is already calculated returned.

```

predFingerFor(i, reqNode,
  targID)  $\triangleq$  predFingerForBody(i, reqNode, targID)
  +
  .predFingerForStop

```

If *break* is false, the return criterion is checked, which may lead to setting *break* to true. predFingerForReturn matches the return criterion, predFingerForContinue continues the loop iteration.

```

predFingerForBody(i,
  reqNode, targID)  $\triangleq$  qry("break", false)@self
  .(predFingerForReturn(i, reqNode, targID)
    +
    predFingerForContinue(i, reqNode, targID))

```

$$\text{predFingerForReturn}(i, \text{reqNode}, \text{targID}) \triangleq \mathbf{qry}(\text{"fingerTable"}, !\text{fingers} \\ \quad [\text{fingers.get}(i).\text{node} \in (\text{self}, \text{targID})])@\text{self} \\ \quad .\mathbf{put}(\text{"break"}, \text{true})@\text{self} \\ \quad .\mathbf{put}(\text{"predFingerResp"}, \text{targID}, \\ \quad \quad \{\text{fingers.get}(i).\text{node}\})@\text{reqNode} \\ \quad .\text{nil}$$

$$\text{predFingerForContinue}(i, \text{reqNode}, \text{targID}) \triangleq \mathbf{qry}(\text{"fingerTable"}, !\text{fingers} \\ \quad [\text{fingers.get}(i).\text{node} \notin (\text{self}, \text{targID})])@\text{self} \\ \quad .\text{nil}$$

$$\text{predFingerForStop} \triangleq \mathbf{qry}(\text{"break"}, \text{true})@\text{self} \\ \quad .\text{nil}$$

$$\text{predFingerReturn}(\text{reqNode}, \text{targID}) \triangleq (\text{predFingerNoReturn} \\ \quad + \\ \quad \text{predFingerReturnSelf}(\text{reqNode}, \text{targID}))$$

After the loop it has to be considered, whether self has to be returned or not.

$$\text{predFingerNoReturn} \triangleq \mathbf{get}(\text{"break"}, \text{true})@\text{self} \\ \quad .\text{predFinger}$$

$$\text{predFingerReturnSelf}(\text{reqNode}, \text{targID}) \triangleq \mathbf{get}(\text{"break"}, \text{false})@\text{self} \\ \quad .\mathbf{put}(\text{"predFingerResp"}, \text{targID}, \text{self})@\text{reqNode} \\ \quad .\text{predFinger}$$

## 4.2 Node Joins

A joining node runs the following procedure:

```

n.join(n')
  if (n')
    init_finger_table(n');
    update_others();
    dataMap = finger[1].node.extract_data(predecessor,n);
  else
    for (i = 1 to m)
      finger[i].node = n
      predecessor = n;

n.init_finger_table(n')
  finger[1].node = n'.find_successor(finger[1].start)
  predecessor = finger[1].node.predecessor;
  finger[1].node.predecessor = n;
  for (i = 1 to 1)
    if (finger[i + 1].start ∈ [n,finger[i].node))

```

```

        finger[i + 1].node = finger[i].node;
    else
        finger[i + 1].node = n'.find_successor(finger[i + 1].start);
n.update_others()
  for (i = 1 to m)
    p = find_predecessor_incl(n - 2i-1);
    p.update_finger_table(n,i);
n.update_finger_table(s,i)
  if (s ∈ [n,finger[i].node))
    finger[i].node = s;
    p = predecessor;
    p.update_finger_table(s,i);

```

We add the *extract\_data* procedure. This procedure transfers all data now belonging to the new node from the successor node. It is called after *update\_others()*.

```

n.extract_data(p,n')
  return dataMap.extract((p,n'));

```

In SCEL, the join procedure is modelled by the process *join(helper)*. Helper is an arbitrary node already participating in the network and is necessary to initialize the fingers list of the new node. *join(helper)* is split into the following tasks:

- *initFingerTable* - initialize finger list via helper.
- *setPred* - set predecessors of self and successor.
- *updateOthers* - update the finger lists of preceding nodes whose fingers could be affected.
- *transferKeys* - get dataMap parts of successor which are now stored at the new node.

Afterwards, the main node process is started.

## join

Join is modelled as follows.

$$\text{join(helper)} \triangleq \text{initFingerTable(helper)} \cdot \text{setPred} \cdot \text{updateOthers} \cdot \text{transferKeys}$$

## init\_finger\_table

Here the whole finger table including the first entry is set. The initialization of the predecessor happens afterwards. This order does not affect the correctness. In *initFingerTable* there is a loop setting the node values of all fingers by calling *findSucc* at the helper node.

$$\text{initFingerTable(helper)} \triangleq \text{put}(\text{"fingerTable"}, \{\text{newList}(m)\})@self \\ \quad \cdot (i : 1 \rightarrow m)(\text{setFinger}(i, \text{fingers}, \text{helper}), \text{nil})$$



```

setFinger(i, helper)  $\triangleq$  get("fingerTable", !fingers)@self
    .put("findSuccReq", self, {self + 2i-1})@helper
    .get("findSuccResp", {self + 2i-1}, fingerNode)@self
    .put("fingerTable", {fingers.set(i, fingerNode)})@self
    .nil

```

setPred sets the attributes *pred* and *succ* and registers the node at the successor as new predecessor.

```

setPred  $\triangleq$  qry("fingerTable", !fingers)@self
    .put("succ", {fingers.get(1).node})@self
    .get("succ", !succ)@self
    .qry("pred", !pred)@succ
    .put("pred", pred)@self
    .put("pred", self)@succ
    .nil

```

### update\_others

updateOthers notifies all preceding nodes, that might have to update their finger table entries. To determine the procedure callees, it uses the process fingerPredIncl.

```

updateOthers  $\triangleq$  (i : 1  $\rightarrow$  m)(updateOther(i), nil)

```

```

updateOther(i)  $\triangleq$  put("findPredInclReq", self, {self - 2i-1})@self
    .get("findPredInclResp", {self - 2i-1}, !targNode)@self
    .put("updateFinger", self, i)@targNode
    .nil

```

### transfer\_keys

Process transferKeys retrieves all data from the successor, that is now managed at this node.

```

transferKeys  $\triangleq$  put("dataMap", {newMap((pred, self)})@self
    .qry("dataMap", !map)@self
    .qry("fingerTable", !fingers)@self
    .put("succ", {fingers.get(1).node})@self
    .get("succ", !succ)@self
    .put("getDataReq", self, {map.domain()})@succ
    .get("getDataResp", {fingers.get(1).node}, !data)@self
    .put("dataMap", {map.addAll(data)})@self
    .nodeProcess

```

### auxiliary processes

To ensure the correct joining behavior, every node has to run the processes `updateFinger` and `getDataHelper`. `updateFinger` is executed at nodes which might need to update their finger list.

$$\begin{aligned} \text{updateFinger} \triangleq & \mathbf{get}(\text{"updateFinger"}, !\text{newID}, !i)@self \\ & .(\text{updateFingerBody}(\text{newID}, i) \\ & + \\ & \text{updateFingerEnd}(\text{newID}, i)) \end{aligned}$$

$$\begin{aligned} \text{updateFingerBody}(\text{newID}, i) \triangleq & \mathbf{get}(\text{"fingerTable"}, !\text{fingers}[\text{newID} \in \\ & [\text{self}, \{\text{fingers.get}(i).\text{node}\}]]@self \\ & .\mathbf{put}(\text{"fingerTable"}, \\ & \{\text{fingers.set}(i, \text{newID})\})@self \\ & .\mathbf{qry}(\text{"pred"}, !\text{pred})@self \\ & .\mathbf{put}(\text{"updateFinger"}, \text{newID}, i)@pred \\ & .\text{updateFinger} \end{aligned}$$

$$\begin{aligned} \text{updateFingerEnd}(\text{newID}, i) \triangleq & \mathbf{qry}(\text{"fingerTable"}, !\text{fingers}[\text{newID} \notin \\ & [\text{self}, \{\text{fingers.get}(i).\text{node}\}]]@self \\ & .\text{updateFinger} \end{aligned}$$

`getDataHelper` provides functionality for extracting data out of `dataMap` to distribute it to new nodes.

$$\begin{aligned} \text{getDataHelper} \triangleq & \mathbf{get}(\text{"getDataReq"}, !\text{reqNode}, \\ & !\text{mapDomain})@self \\ & .\mathbf{qry}(\text{"dataMap"}, !\text{map})@self \\ & .\mathbf{put}(\text{"getDataResp"}, self, \\ & \{\text{map.extract}(\text{mapDomain})\})@reqNode \\ & .\text{getDataHelper} \end{aligned}$$

### 4.3 Leaving nodes

A leaving node executes process `leave` to update the predecessor node of the successor and transmit all data. The pseudocode we defined, looks like this.

```

n.leave()
  finger[1].node.transmit_data(dataMap)
  finger[1].node.predecessor = predecessor;
  for (i = 1 to m)
    p = find_predecessor(n - 2i-1);
    p.clean_finger(n, i, finger[1].node);

n.clean_finger(s, i, c)
  if (finger[i].node == s)

```

```

finger[i].node = c;
p = predecessor;
p.clean_finger(s,i,c);

```

```

n.transmit_data(otherMap)
  dataMap.merge(otherMap);

```

In the following, we describe the SCEL processes.

### leave

Procedure *put\_data* is realized by simply putting the data to the succeeding node. The receiving process is modelled below. Then the preceding nodes are notified about the leave.

$$\text{leave} \triangleq \text{qry}(\text{"fingerTable"}, !\text{fingers})@self$$

```

  .qry("dataMap", !map)@self
  .put("succ", {fingers.get(1).node})@self
  .get("succ", !succ)@self
  .put("predData", map)@succ
  .qry("pred", !pred)@self
  .put("pred", pred)@succ
  .(i : 1 → m)(cleanUpNotify(i, succ), nil)

```

*cleanUpNotify* works similar to *updateOthers* in the joining case.

$$\text{cleanUpNotify}(i, \text{succ}) \triangleq \text{put}(\text{"findPredInclReq"}, self, \{self - 2^{i-1}\})@self$$

```

  .get("findPredInclResp", {self - 2^{i-1}}, !targNode)@self
  .put("cleanUpFinger", self, succ, i)@targNode
  .nil

```

Every node has to run the process *leaveHelper*. It simply waits for incoming data and adds it to the local data map.

$$\text{leaveHelper} \triangleq \text{get}(\text{"predData"}, !\text{predMap})@self$$

```

  .get("dataMap", !map)@self
  .put("dataMap", {map.merge(predMap)})@self
  .leaveHelper

```

### clean\_up\_finger

*cleanUpFinger* is the process receiving the cleanup request. It checks, if the old finger table entry is set to the leaving node. That is decided via a nondeterministic choice. If an update has to take place, the call is propagated backwards

$$\text{cleanUpFinger} \triangleq \text{get}(\text{"cleanUpFinger"}, !\text{leaver}, !\text{cand}, !i)@self$$

```

  .(cleanUpFingerBody(leaver, cand, i)
  +
  cleanUpFingerEnd(leaver, cand, i))

```

$$\text{cleanUpFingerBody}(\text{leaver}, \text{cand}, i) \triangleq \text{get}(\text{"fingerTable"}, !\text{fingers}[\text{leaver} == \text{\{fingers.get(i).node\}}])@self$$

$$\text{.put}(\text{"fingerTable"}, \text{\{fingers.set(i, cand)\}})@self$$

$$\text{.qry}(\text{"pred"}, !\text{pred})@self$$

$$\text{.put}(\text{"cleanUpFinger"}, \text{leaver}, \text{cand}, i)@pred$$

$$\text{.cleanUpFinger}$$

$$\text{cleanUpFingerEnd}(\text{leaver}, \text{cand}, i) \triangleq \text{qry}(\text{"fingerTable"}, !\text{fingers}[\text{leaver}! = \text{\{fingers.get(i).node\}}])@self$$

$$\text{.cleanUpFinger}$$

#### 4.4 Data Lookup

Process `lookupDataClient(key)` is executed for querying the data value stored at `key`. The node responsible is identified and then asked for the data. This section is not part of Chord. We define the following procedure:

```

n.lookup_data(key)
  if (key ∈ dataMap.domain())
    return dataMap.get(key)
  else
    n' = find_successor(key);
    return n'.lookup_data(key);

```

##### lookup\_data client

This is our `lookupDataClient`, which determines the responsible node by calling `findSucc` and then requests the data.

$$\text{lookupDataClient}(\text{key}) \triangleq \text{put}(\text{"findSuccReq"}, \text{self}, \text{key})@self$$

$$\text{.get}(\text{"findSuccResp"}, \text{key}, !\text{targNode})@self$$

$$\text{.put}(\text{"lookupDataReq"}, \text{self}, \text{key})@targNode$$

$$\text{.get}(\text{"lookupDataResp"}, \text{key}, !\text{data})@self$$

$$\text{.nil}$$

##### lookup\_data server

In order to provide data to the network, every node has to run the process `lookupDataServer`. It behaves like a simple getter procedure.

```
lookupDataServer  $\triangleq$  get("lookupDataReq", !reqNode, !key)@self
    .qry("dataMap", !map)@self
    .put("lookupDataResp", key, {map.get(key)})@reqNode
    .lookupDataServer
```

## 4.5 Data Storage

Storing a new key-value pair is achieved through the process `storeDataClient(key, data)`. After identifying the responsible node, the data is transmitted. This section is not part of Chord either. We define the following procedure:

```
n.store_data(key, data)
  if (key  $\in$  dataMap.domain())
    dataMap.put(key, data)
  else
    n' = find_successor(key);
    n'.store_data(key, data);
```

### store\_data client

The SCEL model is split up into a server and a client. The client process looks up the target and then sends the data.

```
storeDataClient(key, data)  $\triangleq$  put("findSuccReq", self, key)@self
    .get("findSuccResp", key, !targNode)@self
    .put("storeDataReq", key, data)@targNode
    .nil
```

### lookup\_data server

Every node has to run a `storeDataServer` process, waiting for new key-value pairs to arrive.

```
storeDataServer  $\triangleq$  get("storeDataReq", !key, !data)@self
    .qry("dataMap", !map)@self
    .put("dataMap", {map.put(key, data)})@self
    .storeDataServer
```

## 4.6 Stabilization

To react to node failures, the finger tables have to be checked periodically. The processes `stabilize` and `fix_fingers` update the successor and the other fingers.

```

n.stabilize()
  x = finger[1].node.predecessor;
  if ( $x \in (n, \text{successor})$ )
    finger[1].node = x;
  finger[1].node.notify(self);

n.notify(n')
  if (predecessor is nil or  $n' \in (\text{predecessor}, n)$ )
    predecessor = n';

n.fix_fingers()
  i = random index > 1 into finger[];
  finger[i].node = find_successor(finger[i].start);

```

These procedures are expressed in SCEL as follows:

### stabilize

A stabilizing node checks if the predecessor of the successor is itself. If not, the new node in between is the new successor and is notified.

$$\text{stabilize} \triangleq \mathbf{qry}(\text{"fingerTable"}, !\text{fingers})@self$$

$$\quad .\mathbf{put}(\text{"succ"}, \{\text{fingers.get}(1).\text{node}\})@self$$

$$\quad .\mathbf{get}(\text{"succ"}, !\text{succ})@self$$

$$\quad .(\text{newSuccNode}(\text{succ}) + \text{notifySucc}(\text{succ}))$$

$$\text{newSuccNode}(\text{succ}) \triangleq \mathbf{qry}(\text{"pred"}, !x[x \in (\text{self}, \text{succ})])@succ$$

$$\quad .\mathbf{get}(\text{"fingerTable"}, !\text{fingers})@self$$

$$\quad .\mathbf{put}(\text{"fingerTable"}, \{\text{fingers.set}(1, x)\})@self$$

$$\quad .\mathbf{put}(\text{"predNotify"}, \text{self})@x$$

$$\quad .\text{stabilize}$$

$$\text{notifySucc}(\text{succ}) \triangleq \mathbf{qry}(\text{"pred"}, !x[x \notin (\text{self}, \text{succ})])@succ$$

$$\quad .\mathbf{put}(\text{"predNotify"}, \text{self})@succ$$

$$\quad .\text{stabilize}$$

If a node is notified about a potentially new predecessor is realized by predNotify.

$$\text{predNotify} \triangleq \mathbf{get}(\text{"predNotify"}, !\text{newPred})@self$$

$$\quad .(\text{predChange}(\text{newPred}) + \text{predNoChange}(\text{newPred}))$$

$$\text{predChange}(\text{newPred}) \triangleq \mathbf{qry}(\text{"pred"}, !\text{pred}[\text{pred} == \text{nil} \parallel \text{newPred} \in (\text{pred}, \text{self})])@self$$

$$\quad .\mathbf{put}(\text{"pred"}, \text{newPred})@self$$

$$\quad .\text{predNotify}$$

$$\text{predNoChange}(\text{newPred}) \triangleq \mathbf{qry}(\text{"pred"}, !\text{pred}[\text{pred!} = \text{nil} \ \&\& \ \text{newPred} \notin (\text{pred}, \text{self})])@\text{self} \\ \text{.predNotify}$$

### fix\_fingers

Process `fixFinger` checks if the node value of the finger table entry at given index is correct by calculating the successor of the start value again.

$$\text{fixFinger}(\text{index}) \triangleq \mathbf{get}(\text{"fingerTable"}, !\text{fingers})@\text{self} \\ \text{.put}(\text{"findSuccReq"}, \{\text{self} + 2^{(i-1)}\}, \text{self})@\text{self} \\ \text{.get}(\text{"findSuccResp"}, \{\text{self} + 2^{(i-1)}\}, !\text{succ}) \\ \text{.put}(\text{"fingerTable"}, \{\text{fingers.set}(\text{index}, \text{succ})\})@\text{self}$$

## 4.7 Main Node Process

This process contains all auxiliary processes to ensure the Chord functionality.

$$\text{nodeProcess} \triangleq \text{findSucc} \mid \text{findPredExcl} \mid \text{findPredIncl} \\ \mid \text{getSucc} \mid \text{predFinger} \mid \text{updateFinger} \\ \mid \text{getDataHelper} \mid \text{leaveHelper} \mid \text{lookupDataServer} \\ \mid \text{storeDataServer} \mid \text{cleanUpFinger}$$

Additionally, the following processes are called at some point by external means (timer-based or by user action).

- `storeDataClient(key,data)`
- `lookupDataClient(key)`
- `leave`
- `join(helper)`
- `stabilize`
- `fixFinger(index)`

## 5 Implementing Chord

When implementing the Chord protocol in Java, we focussed on keeping the program as close as possible to our SCEL model. In this section we explain the implementation starting with a short description of the chosen architecture. We finish with a short usage documentation and screenshots.

### 5.1 Architecture

Chord is implemented in Java using the standard Java API. No external libraries or frameworks are used to ensure that the code is easily understandable. The basic architecture is three-layered and shown in figure 1. The classes shown in the layers represent a view into the system, they are supported by others (see source code for a detailed view). The layers are discussed in the following subsections.

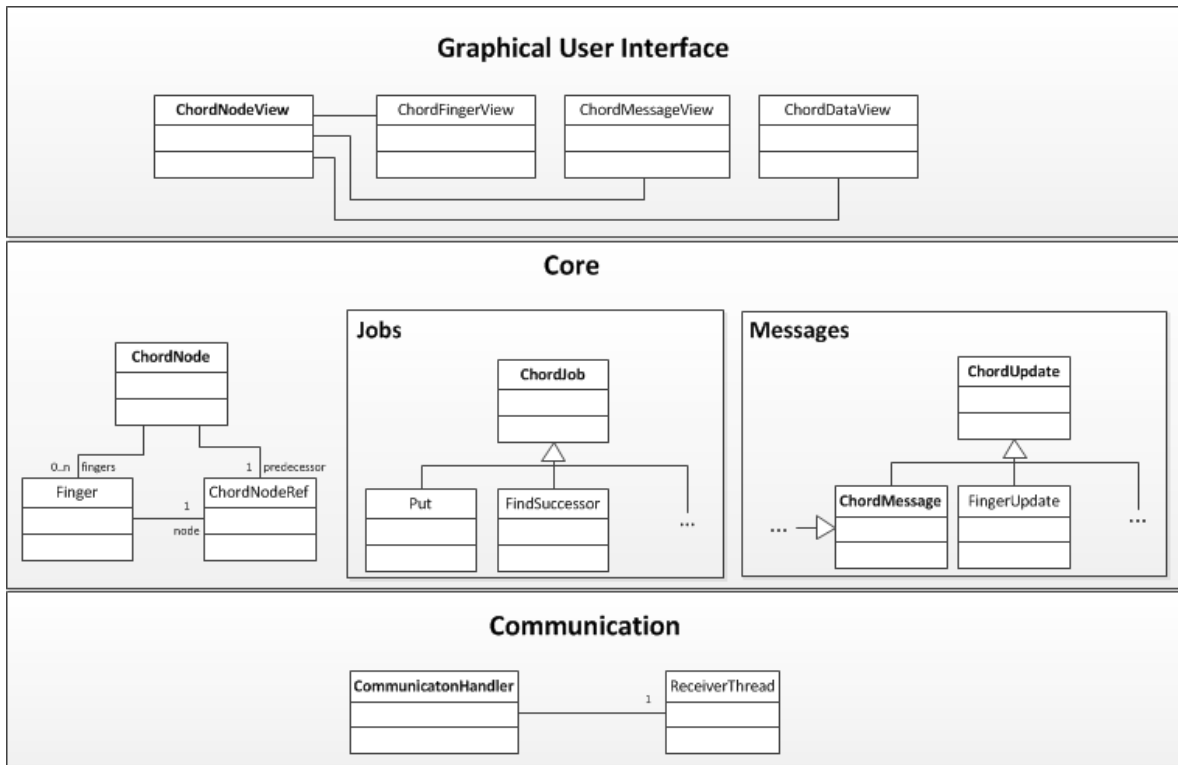


Figure 1: Layered Architecture

### 5.1.1 Core layer

The main part of the implementation is the *core*, which forms the Chord layer. It contains all the protocol logic. A node is represented by a *ChordNode* instance. A *ChordNode* object has links to its fingers and its predecessor. To the outside, it provides the procedures *put* and *get*, which allow adding a key value pair to the ring or removing it.

The procedures each node has to execute from time to time are *ChordJob* objects. Each job is a thread which is started individually on demand, and runs concurrent to all other jobs.

Communication at this layer is realized through serializable *ChordMessage* objects. Each message type is represented by its own class, which may contain additional information required for the receiver.

Remote procedure calls are also realized through messages. If a procedure has to be run on another node, a *ChordMessage* is used to trigger the required behavior. Procedures with result values can be accessed via queries.

If a job wants to start a query and awaits remotely calculated answers, it uses *ChordQuery* objects. Matching of query responses to waiting jobs happens at this layer. A node is identified by its Chord id at the core layer and by an ip address and a port number at the communication layer. *ChordNodeRef* objects which represent network identification contain these three values.

### 5.1.2 Communication layer

The layer underneath *core* is the communication layer. It provides functionality for sending UDP messages to a given receiver address. The main part in this layer is provided by *CommunicationHandlers*.



The class *CommunicationHandler* defines the method *sendMessage*, which serializes a given object and sends it to a specified address. Incoming messages are received by a dedicated *ReceiverThread*, which notifies about new packets via the observer pattern.

### 5.1.3 Graphical User Interface

The Chord implementation contains a graphical user interface for interacting with the implementation.

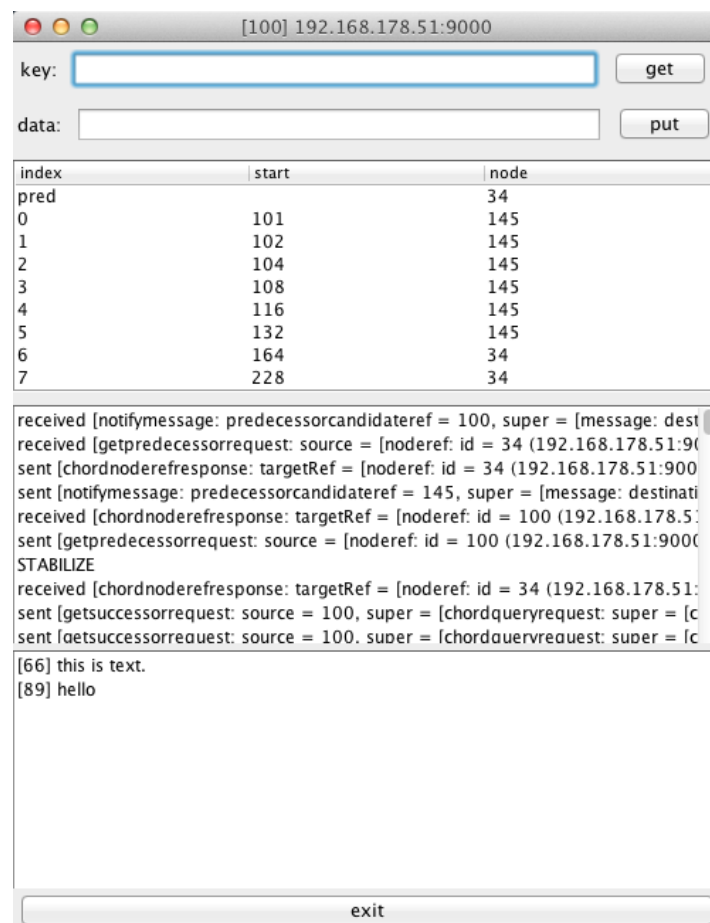


Figure 2: Screenshot of a Chord Node

Figure 2 shows the GUI of Chord nodes. The first section with textfields and buttons allows to get and put data. The string in the key textfield will be hashed to a key.

The table below shows current predecessor and finger nodes of the node. Each finger is displayed with its start value and the actual node which is addressed here.

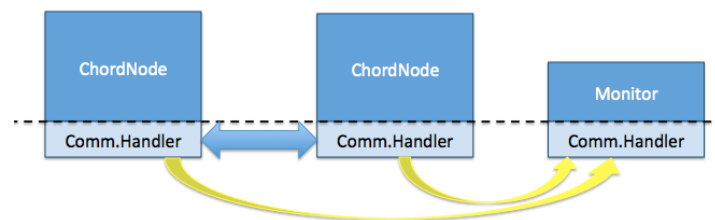
The two lists below show all received and sent message as well as some updates (example: 'STABILIZE') and all data that is stored at this node in the form [key] data.

The exit button on the bottom of the screen allows controlled leaves. Pressing it starts the Chord job *leave* and then terminates the program.

## 5.2 Monitoring the System

To be able to figure out what is going inside the Chord network from a global perspective, the Chord implementation contains an additional *monitoring* component. Although the main Chord protocol is strictly peer-to-peer, we add one global component purely for monitoring which receives all updates sent within the network as well, and is able to display them graphically.

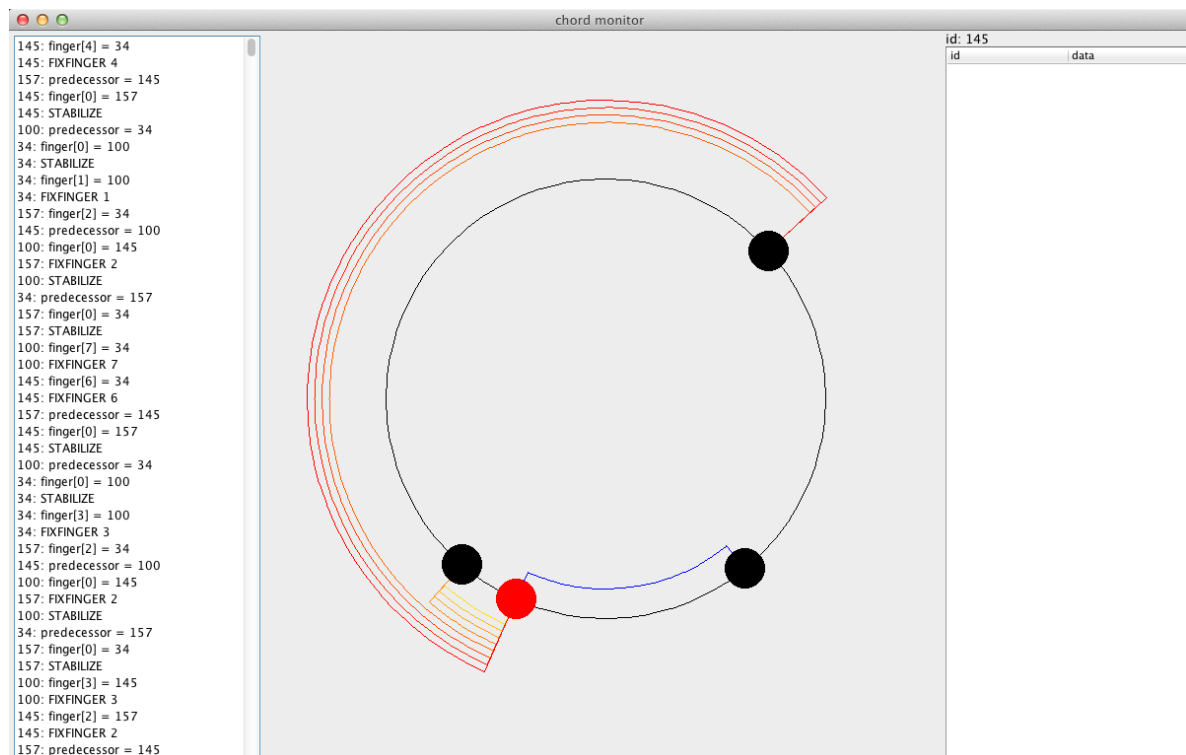
A Chord monitor runs independently on a machine in the network at a certain address and port. It receives *ChordUpdates*, which are sent by the nodes at relevant events like joins, leaves, data change-ment, stabilizes and more.



A monitor is a passive component. It only listens for updates and displays them. If data information is incomplete, the monitor does not fetch data on its own. The flow of data in the Chord implementation is shown in figure 5.2.

An example of the monitor UI in action is shown in figure 5.2. The middle section shows an image of the logical ring. Here, four rings are participating. Now, node 145 is activated. The node's predecessor is shown in blue; its fingers are displayed in yellow to red around the ring. On the right, the id and stored data (currently none) are displayed.

On the left of the window, a list of incoming *ChordUpdates* can be viewed. These updates actually provide all information the monitor has, as stated above.



## 6 Conclusion

This document has shown how Chord, a peer-to-peer network lookup protocol, can be modelled in the SCEL modelling language, and implemented in plain Java.

During this work, we have identified a shortcoming in the Chord protocol as discussed in section 2. We have furthermore discussed our SCEL dialect and added some macros which might be beneficial in the future for specifying other service components (section 3).

Our model of the Chord protocol in SCEL can be found in section 4, which also includes some additional Chord procedures we required for fully expressing a system based on Chord (for example, for leaving nodes). Our implementation has been discussed in section 5, where we have also described a monitor for simplifying the visualization of the Chord ring.

## References

- [dNFLP11] Rocco de Nicola, Gianluigi Ferrari, Michele Loreti, and Roberto Pugliese. D1.1: Language Primitives for Coordination, Resource Negotiation, and Task Description. ASCENS Deliverable, November 2011.
- [GLPT12] Edmond Gjondrekaj, Michele Loreti, Rosario Pugliese, and Francesco Tiezzi. Modeling adaptation with a tuple-based coordination language. In *Proceedings of the 27th Annual ACM Symposium on Applied Computing, SAC '12*, pages 1522–1527, New York, NY, USA, 2012. ACM.
- [SMLN<sup>+</sup>03] Ion Stoica, Robert Morris, David Liben-Nowell, David R. Karger, M. Frans Kaashoek, Frank Dabek, and Hari Balakrishnan. Chord: a scalable peer-to-peer lookup protocol for internet applications. *IEEE/ACM Trans. Netw.*, 11(1):17–32, February 2003.
- [vRA<sup>+</sup>11] Nikola Šerbedžija, Stephan Reiter, Maximilian Ahrens, José Velasco, Carlo Pinciroli, Nicklas Hoch, and Bernd Werther. D7.1: First Report on WP7: Requirement Specification and Scenario Description of the ASCENS Case Studies. ASCENS Deliverable, November 2011.