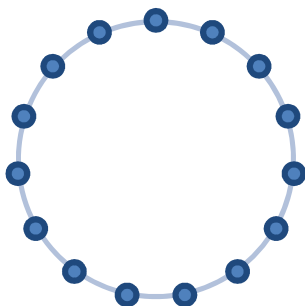


## TRANSACTIONS



## Scalaris:

### Users and Developers Guide

Version 0.3.0+svn

December 13, 2011

Copyright 2007-2011 Zuse Institute Berlin and onScale solutions.

Licensed under the Apache License, Version 2.0 (the "License"); you may not use this file except in compliance with the License. You may obtain a copy of the License at

<http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

# Contents

<b>I. Users Guide</b>	<b>5</b>
<b>1. Introduction</b>	<b>6</b>
1.1. Brewer's CAP Theorem . . . . .	6
1.2. Scientific Background . . . . .	7
<b>2. Download and Installation</b>	<b>8</b>
2.1. Requirements . . . . .	8
2.2. Download . . . . .	8
2.2.1. Development Branch . . . . .	8
2.2.2. Releases . . . . .	8
2.3. Build . . . . .	9
2.3.1. Linux . . . . .	9
2.3.2. Windows . . . . .	9
2.3.3. Java-API . . . . .	9
2.3.4. Python-API . . . . .	10
2.3.5. Ruby-API . . . . .	10
2.4. Installation . . . . .	10
<b>3. Setting up Scalaris</b>	<b>12</b>
3.1. Runtime Configuration . . . . .	12
3.1.1. Logging . . . . .	12
3.2. Running Scalaris . . . . .	13
3.2.1. Running on a local machine . . . . .	13
3.2.2. Running distributed . . . . .	14
3.3. Custom startup using scalarisctl . . . . .	14
<b>4. Using the system</b>	<b>15</b>
4.1. Application Programming Interfaces (APIs) . . . . .	15
4.1.1. Supported Types . . . . .	16
4.1.2. JSON API . . . . .	17
4.1.3. Java API . . . . .	20
4.2. Command Line Interfaces . . . . .	20
4.2.1. Java command line interface . . . . .	20
4.2.2. Python command line interface . . . . .	21
4.2.3. Ruby command line interface . . . . .	22
<b>5. Testing the system</b>	<b>23</b>
5.1. Erlang unit tests . . . . .	23
5.2. Java unit tests . . . . .	23
5.3. Python unit tests . . . . .	24
5.4. Interoperability Tests . . . . .	26

<b>6. Troubleshooting</b>	<b>27</b>
6.1. Network . . . . .	27
6.2. Miscellaneous . . . . .	27
 <b>II. Developers Guide</b>	 <b>28</b>
<b>7. General Hints</b>	<b>29</b>
7.1. Coding Guidelines . . . . .	29
7.2. Testing Your Modifications and Extensions . . . . .	29
7.3. Help with Digging into the System . . . . .	29
 <b>8. System Infrastructure</b>	 <b>30</b>
8.1. Groups of Processes . . . . .	30
8.2. The Communication Layer <code>comm</code> . . . . .	30
8.3. The <code>gen_component</code> . . . . .	30
8.3.1. A basic <code>gen_component</code> including a message handler . . . . .	31
8.3.2. How to start a <code>gen_component</code> ? . . . . .	32
8.3.3. When does a <code>gen_component</code> terminate? . . . . .	32
8.3.4. What happens when unexpected events / messages arrive? . . . . .	33
8.3.5. What if my message handler generates an exception or crashes the process? . . . . .	33
8.3.6. Changing message handlers and implementing state dependent message re- sponsiveness as a state-machine . . . . .	33
8.3.7. Handling several messages atomically . . . . .	34
8.3.8. Halting and pausing a <code>gen_component</code> . . . . .	34
8.3.9. Integration with <code>pid_groups</code> : Redirecting messages to other <code>gen_components</code> . . . . .	34
8.3.10. Replying to ping messages . . . . .	35
8.3.11. The debugging interface of <code>gen_component</code> : Breakpoints and step-wise exe- cution . . . . .	35
8.3.12. Future use and planned extensions for <code>gen_component</code> . . . . .	38
8.4. The Process' Database ( <code>pdb</code> ) . . . . .	38
8.5. Failure Detectors ( <code>fd</code> ) . . . . .	38
8.6. Monitoring Statistics ( <code>monitor</code> , <code>rrd</code> ) . . . . .	38
8.7. Writing Unittests . . . . .	39
8.7.1. Plain unittests . . . . .	39
8.7.2. Randomized Testing using <code>tester.erl</code> . . . . .	39
 <b>9. Basic Structured Overlay</b>	 <b>40</b>
9.1. Ring Maintenance . . . . .	40
9.2. T-Man . . . . .	40
9.3. Routing Tables . . . . .	40
9.3.1. The routing table process ( <code>rt_loop</code> ) . . . . .	43
9.3.2. Simple routing table ( <code>rt_simple</code> ) . . . . .	43
9.3.3. Chord routing table ( <code>rt_chord</code> ) . . . . .	46
9.4. Local Datastore . . . . .	50
9.5. Cyclon . . . . .	50
9.6. Vivaldi Coordinates . . . . .	50
9.7. Estimated Global Information (Gossiping) . . . . .	50
9.8. Load Balancing . . . . .	50
9.9. Broadcast Trees . . . . .	50

<b>10. Transactions in Scalaris</b>	<b>51</b>
10.1. The Paxos Module . . . . .	51
10.2. Transactions using Paxos Commit . . . . .	51
10.3. Applying the Tx-Modules to replicated DHTs . . . . .	51
<b>11. How a node joins the system</b>	<b>52</b>
11.1. Supervisor-tree of a Scalaris node . . . . .	53
11.2. Starting the sup_dht_node supervisor and general processes of a node . . . . .	53
11.3. Starting the sup_dht_node_core supervisor with a peer and some paxos processes . .	55
11.4. Initializing a dht_node-process . . . . .	56
11.5. Actually joining the ring . . . . .	56
11.5.1. A single node joining an empty ring . . . . .	56
11.5.2. A single node joining an existing (non-empty) ring . . . . .	57
<b>12. Directory Structure of the Source Code</b>	<b>65</b>
<b>13. Java API</b>	<b>66</b>

Part I.

# Users Guide

# 1. Introduction

Scalaris is a scalable, transactional, distributed key-value store based on the peer-to-peer principle. It can be used to build scalable Web 2.0 services. The concept of Scalaris is quite simple: Its architecture consists of three layers.

It provides self-management and scalability by replicating services and data among peers. Without system interruption it scales from a few PCs to thousands of servers. Servers can be added or removed on the fly without any service downtime.



Scalaris takes care of:

- Fail-over
- Data distribution
- Replication
- Strong consistency
- Transactions

The Scalaris project was initiated by Zuse Institute Berlin and onScale solutions and was partly funded by the EU projects Selfman and XtreamOS. Additional information (papers, videos) can be found at <http://www.zib.de/CSR/Projects/scalaris> and <http://www.onscale.de/scalarix.html>.

## 1.1. Brewer's CAP Theorem

In distributed computing there exists the so called CAP theorem. It basically says that there are three desirable properties for distributed systems but one can only have any two of them.

**Strict Consistency.** Any read operation has to return the result of the latest write operation on the same data item.

**Availability.** Items can be read and modified at any time.

**Partition Tolerance.** The network on which the service is running may split into several partitions which cannot communicate with each other. Later on the networks may re-join again.

For example, a service is hosted on one machine in Seattle and one machine in Berlin. This service is partition tolerant if it can tolerate that all Internet connections over the Atlantic (and Pacific) are interrupted for a few hours and then get repaired.

The goal of Scalaris is to provide strict consistency and partition tolerance. We are willing to sacrifice availability to make sure that the stored data is always consistent. I.e. when you are running Scalaris with a replication degree of 4 and the network splits into two partitions, one partition with three replicas and one partition with one replica, you will be able to continue to use the service only in the larger partition. All requests in the smaller partition will time out until the two networks merge again. Note, most other key-value stores tend to sacrifice consistency.

## 1.2. Scientific Background

**Basics.** The general structure of Scalaris is modelled after Chord. The Chord paper [4] describes the ring structure, the routing algorithms, and basic ring maintenance.

The main routines of our Chord node are in `src/dht_node.erl` and the join protocol is implemented in `src/dht_node_join.erl` (see also Chap. 11 on page 52). Our implementation of the routing algorithms is described in more detail in Sect. 9.3 on page 40 and the actual implementation is in `src/rt_chord.erl`.

**Transactions.** The most interesting part is probably the transaction algorithms. The most current description of the algorithms and background is in [6].

The implementation consists of the paxos algorithm in `src/paxos` and the transaction algorithms itself in `src/transactions` (see also Chap. 10 on page 51).

**Ring Maintenance.** We changed the ring maintenance algorithm in Scalaris. It is not the standard Chord one, but a variation of T-Man [5]. It is supposed to fix the ring structure faster. In some situations, the standard Chord algorithm is not able to fix the ring structure while T-Man can still fix it. For node sampling, our implementation relies on Cyclon [7].

The T-Man implementation can be found in `src/rm_tman.erl` and the Cyclon implementation in `src/cyclon`.

**Vivaldi Coordinates.** For some experiments, we implemented so called Vivaldi coordinates [2]. They can be used to estimate the network latency between arbitrary nodes.

The implementation can be found in `src/vivaldi.erl`.

**Gossiping.** For some algorithms, we use estimates of global information. These estimates are aggregated with the help of gossiping techniques [8].

The implementation can be found in `src/gossip.erl`.

## 2. Download and Installation

### 2.1. Requirements

For building and running Scalaris, some third-party software is required which is not included in the Scalaris sources:

- Erlang R13B01 or newer
- OpenSSL (required by Erlang's crypto module)
- GNU-like Make and autoconf (not required on Windows)

To build the Java API (and its command-line client) the following programs are also required:

- Java Development Kit 6
- Apache Ant

Before building the Java API, make sure that JAVA\_HOME and ANT\_HOME are set. JAVA\_HOME has to point to a JDK installation, and ANT\_HOME has to point to an Ant installation.

To build the Python API (and its command-line client) the following programs are also required:

- Python  $\geq$  2.6

### 2.2. Download

The sources can be obtained from <http://code.google.com/p/scalaris>. RPM and DEB packages are available from <http://download.opensuse.org/repositories/home:/scalaris/> for various Linux distributions.

#### 2.2.1. Development Branch

You find the latest development version in the svn repository:

```
# Non-members may check out a read-only working copy anonymously over HTTP.  
svn checkout http://scalaris.googlecode.com/svn/trunk/ scalaris-read-only
```

#### 2.2.2. Releases

Releases can be found under the 'Download' tab on the web-page.



## 2.3. Build

### 2.3.1. Linux

Scalaris uses autoconf for configuring the build environment and GNU Make for building the code.

```
%> ./configure
%> make
%> make docs
```

For more details read README in the main Scalaris checkout directory.

### 2.3.2. Windows

We are currently not supporting Scalaris on Windows. However, we have two small .bat files for building and running Scalaris nodes. It seems to work but we make no guarantees.

- Install Erlang  
<http://www.erlang.org/download.html>
- Install OpenSSL (for crypto module)  
<http://www.slproweb.com/products/Win32OpenSSL.html>
- Checkout Scalaris code from SVN
- adapt the path to your Erlang installation in build.bat
- start a cmd.exe
- go to the Scalaris directory
- run build.bat in the cmd window
- check that there were no errors during the compilation; warnings are fine
- go to the bin sub-directory
- adapt the path to your Erlang installation in firstnode.bat, joining\_node.bat
- run firstnode.bat or one of the other start scripts in the cmd window

build.bat will generate a Emakefile if there is none yet. If you have Erlang < R13B04, you will need to adapt the Emakefile. There will be empty lines in the first three blocks ending with “}”.: add the following to these lines and try to compile again. It should work now.

```
, {d, type_forward_declarations_are_not_allowed}
, {d, forward_or_recursive_types_are_not_allowed}
```

For the most recent description please see the FAQ at <http://code.google.com/p/scalaris/wiki/FAQ>.

### 2.3.3. Java-API

The following commands will build the Java API for Scalaris:

```
%> make java
```

This will build `scalaris.jar`, which is the library for accessing the overlay network. Optionally, the documentation can be build:

```
%> cd java-api
%> ant doc
```

### 2.3.4. Python-API

The Python API for Python 2.\* (at least 2.6) is located in the `python-api` directory. Files for Python 3.\* can be created using `2to3` from the files in `python-api`. The following command will use `2to3` to convert the modules and place them in `python3-api`.

```
%> make python3
```

Both versions of python will compile required modules on demand when executing the scripts for the first time. However, pre-compiled modules can be created with:

```
%> make python
%> make python3
```

### 2.3.5. Ruby-API

The Ruby API for Ruby `>= 1.8` is located in the `ruby-api` directory. Compilation is not necessary.

## 2.4. Installation

For simple tests, you do not need to install Scalaris. You can run it directly from the source directory. Note: `make install` will install Scalaris into `/usr/local` and place `scalarisctl` into `/usr/local/bin`, by default. But it is more convenient to build an RPM and install it. On openSUSE, for example, do the following:

```
export SCALARIS_SVN=http://scalaris.googlecode.com/svn/trunk
for package in main bindings; do
  mkdir -p ${package}
  cd ${package}
  svn export ${SCALARIS_SVN}/contrib/packages/${package}/checkout.sh
  ./checkout.sh
  cp * /usr/src/packages/SOURCES/
  rpmbuild -ba scalaris*.spec
  cd ..
done
```

If any additional packages are required in order to build an RPM, `rpmbuild` will print an error.

Your source and binary RPMs will be generated in `/usr/src/packages/SRPMS` and `RPMS`. We build RPM and DEB packages using the latest tagged version as well as checkouts from `svn` and provide them using the Open Build Service at <http://download.opensuse.org/repositories/home:/scalaris/>. Packages are available for

- Fedora 14, 15
- Mandriva 2010, 2010.1,

- openSUSE 11.3, 11.4, Factory, Tumbleweed
- SLE 10, 11, 11SP1,
- CentOS 5.5,
- RHEL 5.5, 6,
- Debian 5.0, 6.0 and
- Ubuntu 10.04, 10.10, 11.04.

An up-to-date list of available repositories can be found at [https://code.google.com/p/scalaris/wiki/FAQ#Prebuild\\_packages](https://code.google.com/p/scalaris/wiki/FAQ#Prebuild_packages).

For those distributions which provide a recent-enough Erlang version, we build the packages using their Erlang package and recommend using the same version that came with the distribution. In this case we do not provide Erlang packages in our repository.

Exceptions are made for openSUSE-based and RHEL-based distributions as well as Debian 5.0:

- For openSUSE, we provide the package from the `devel:languages:erlang` repository.
- For RHEL-based distributions (CentOS 5, RHEL 5, RHEL 6) we included the Erlang package from the EPEL repository of RHEL 6.
- For Debian 5.0 we included the Erlang package of Ubuntu 11.04.

## 3. Setting up Scalaris

*Description is based on SVN revision r1810.*

### 3.1. Runtime Configuration

Scalaris reads two configuration files from the working directory: `bin/scalaris.cfg` (mandatory) and `bin/scalaris.local.cfg` (optional). The former defines default settings and is included in the release. The latter can be created by the user to alter settings. A sample file is provided as `bin/scalaris.local.cfg.example`. To run Scalaris distributed over several nodes, each node requires a `bin/scalaris.local.cfg`:

File `scalaris.local.cfg`:

```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
% Settings for distributed Erlang
% (see scalaris.hrl to switch)

% {mgmt_server, {mgmt_server, 'mgmt_server@foo.bar.com'}}.
% {known_hosts, [{service_per_vm, 'firstnode@foo.bar.com'}]}.

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
% Settings for TCP mode.
% (see scalaris.hrl to switch)

%% userdevguide-begin local_cfg:distributed
% Insert the appropriate IP-addresses for your setup
% as comma separated integers:
% IP Address, Port, and label of the boot server
{mgmt_server, {{127,0,0,1}, 14194, mgmt_server}}.

% IP Address, Port, and label of a node which is already in the system
{known_hosts, [{127,0,0,1}, 14195, service_per_vm]}.
%% userdevguide-end local_cfg:distributed
```

A Scalaris deployment can have a management server and several nodes. The management-server is optional and provides a global view on all nodes of a Scalaris deployment which contact this server, i.e. have its address specified in the `mgmt_server` configuration setting.

In this example, the `mgmt_server`'s location is defined as an IP address plus a TCP port and its Erlang-internal process name. If the deployment should not use a management server, replace the setting with an invalid address, e.g. `' null '`.

#### 3.1.1. Logging

Scalaris uses the `log4erl` library (see `contrib/log4erl`) for logging status information and error messages. The log level can be configured in `bin/scalaris.cfg` for both the stdout and file logger. The default value is `warn`; only warnings, errors and severe problems are logged.

```
%% @doc Loglevel: debug < info < warn < error < fatal < none
```

```
{log_level, warn}.
{log_level_file, warn}.
```

In some cases, it might be necessary to get more complete logging information, e.g. for debugging. In Chapter 11 on page 52, we are explaining the startup process of Scalaris nodes in more detail, here the `info` level provides more detailed information.

```
%% @doc Loglevel: debug < info < warn < error < fatal < none
{log_level, info}.
{log_level_file, info}.
```

## 3.2. Running Scalaris

As mentioned above, Scalaris consists of:

- management servers and
- regular nodes

The management server will maintain a list of nodes participating in the system. A regular node is either the first node in a system or joins an existing system deployment.

### 3.2.1. Running on a local machine

Open at least two shells. In the first, inside the Scalaris directory, start the first node (`firstnode.bat` on Windows):

```
%> ./bin/firstnode.sh
```

This will start a new Scalaris deployment with a single node, including a management server. On success <http://localhost:8000> should point to the management interface page of the management server. The main page will show you the number of nodes currently in the system. A first Scalaris node should have started and the number should show 1 node. The main page will also allow you to store and retrieve key-value pairs but should not be used by applications to access Scalaris. See Section 4.1 on page 15 for application APIs.

In a second shell, you can now start a second Scalaris node. This will be a ‘regular node’:

```
%> ./bin/joining_node.sh
```

The second node will read the configuration file and use this information to contact a number of known nodes (set by the `known_hosts` configuration setting) and join the ring. It will also register itself with the management server. The number of nodes on the web page should have increased to two by now.

Optionally, a third and fourth node can be started on the same machine. In a third shell:

```
%> ./bin/joining_node.sh 2
```

In a fourth shell:

```
%> ./bin/joining_node.sh 3
```

This will add two further nodes to the deployment. The `./bin/joining_node.sh` script accepts a number as its parameter which will be added to the started node's name, i.e. 1 will lead to a node named `node1`. The web pages at <http://localhost:8000> should show the additional nodes.

### 3.2.2. Running distributed

Scalaris can be installed on other machines in the same way as described in Section 2.4 on page 10. In the default configuration, nodes will look for the management server on 127.0.0.1 on port 14195. You should create a `scalaris.local.cfg` pointing to the node running the management server. You should also add a list of known nodes.

File `scalaris.local.cfg`:

```
12 % Insert the appropriate IP-addresses for your setup
13 % as comma separated integers:
14 % IP Address, Port, and label of the boot server
15 {mgmt_server, {{127,0,0,1}, 14194, mgmt_server}}.
16
17 % IP Address, Port, and label of a node which is already in the system
18 {known_hosts, {{127,0,0,1}, 14195, service_per_vm}}}.
```

If you are starting the management server using `firstnode.sh`, it will listen on port 14195 and you have to change the port and the IP address in the configuration file. Otherwise the other nodes will not find the management server. Calling `./bin/joining_node.sh` on a remote machine will start the node and automatically contact the configured management server.

## 3.3. Custom startup using `scalarisctl`

On linux you can also use the `scalarisctl` script to start a management server and 'regular' nodes directly.

```
%> ./bin/sclarisctl -h
```

```
usage: scalarisctl [options] [services] <cmd>
options:
  -h          - print this help message
  -d          - daemonize
  -f          - first node (to start a new Scalaris instead of joining one) (not with -q)
  -q          - elect first node from known hosts (not with -f)
  -n <name>   - Erlang process name (default 'node')
  -p <port>   - TCP port for the Scalaris node
  -y <port>   - TCP port for the built-in webserver
  -k <key>    - join at the given key
  -v          - verbose
services:
  -m          - global Scalaris management server
  -s          - Scalaris node (see also -f)
commands:
  checkinstallation
    - test installation
  start       - start services (see -m and -s)
  stop        - stop a scalaris process defined by its name (see -n)
  restart     - restart a scalaris process by its name (see -n)

  list        - list locally running Erlang VMs
  debug       - connect to a running node via an Erlang shell
```

## 4. Using the system

*Description is based on SVN revision r1936.*

Scalaris can be used with one of the provided command line interfaces or by using one of the APIs in a custom program. The following sections will describe the APIs in general, each API in more detail and the use of our command line interfaces.

### 4.1. Application Programming Interfaces (APIs)

Currently we offer the following APIs:

- an *Erlang API* running on the node Scalaris is run  
(functions can be called using remote connections with distributed Erlang)
- a *Java API* using Erlang's JInterface library  
(connections are established using distributed Erlang)
- a generic *JSON API*  
(offered by an integrated HTTP server running on each Scalaris node)
- a *Python API* for Python  $\geq 2.6$  using JSON to talk to Scalaris.
- a *Ruby API* for Ruby  $\geq 1.8$  using JSON to talk to Scalaris.

Each API contains methods for accessing functions from the three layers Scalaris is composed of. Table 4.1 shows the modules and classes of Erlang, Java, Python and Ruby and their mapping to these layers. The appropriate JSON calls are shown in Section 4.1.2 on page 17.

Special care needs to be taken when trying to delete keys (no matter which API is used). This can only be done outside the transaction layer and is thus not absolutely safe. Refer to the following thread on the mailing list: [http://groups.google.com/group/scalaris/browse\\_thread/thread/ff1d9237e218799](http://groups.google.com/group/scalaris/browse_thread/thread/ff1d9237e218799).

	Erlang module	Java class in de.zib.scalaris	Python / Ruby class in module scalaris
Transaction Layer	api_tx	Transaction, TransactionSingleOp	Transaction, TransactionSingleOp
	api_pubsub	PubSub	PubSub
Replication Layer	api_rdht	ReplicatedDHT	ReplicatedDHT
P2P Layer	api_dht		
	api_dht_raw		
	api_vm	ScalarisVM	

Table 4.1.: Layered API structure

	Erlang	Java	JSON	Python	Ruby
boolean	<code>boolean()</code>	<code>bool</code> , <code>Boolean</code>	<code>true</code> , <code>false</code>	<code>True</code> , <code>False</code>	<code>true</code> , <code>false</code>
integer	<code>integer()</code>	<code>int</code> , <code>Integer</code> <code>long</code> , <code>Long</code> <code>BigInteger</code>	<code>int</code>	<code>int</code>	<code>Fixnum</code> , <code>Bignum</code>
float	<code>float()</code>	<code>double</code> , <code>Double</code>	<code>int frac</code> <code>int exp</code> <code>int frac exp</code>	<code>float</code>	<code>Float</code>
string	<code>string()</code>	<code>String</code>	<code>string</code>	<code>str</code>	<code>String</code>
binary	<code>binary()</code>	<code>byte[]</code>	<code>string</code> (base64-encoded)	<code>bytearray</code>	<code>String</code>
list(type)	<code>[type()]</code>	<code>List&lt;Object&gt;</code>	<code>array</code>	<code>list</code>	<code>Array</code>
JSON	<code>json_obj()</code> *	<code>Map&lt;String, Object&gt;</code>	<code>object</code>	<code>dict</code>	<code>Hash</code>
custom	<code>any()</code>	<code>OtpErlangObject</code>	<code>/</code>	<code>/</code>	<code>/</code>

\*  
`json_obj() :: {struct, [Key::atom() | string(), Value::json_val()]}`  
`json_val() :: string() | number() | json_obj() | {array, [any()]} | true | false | null`

Table 4.2.: Types supported by the Sclaris APIs

#### 4.1.1. Supported Types

Different programming languages have different types. In order for our APIs to be compatible with each other, only a subset of the available types is officially supported.

Keys are always strings. In order to avoid problems with different encodings on different systems, we suggest to only use ASCII characters.

For *values* we distinguish between *native*, *composite* and *custom* types.

*Native* types are

- boolean values
- integer numbers
- floating point numbers
- strings and
- binary objects (a number of bytes).

*Composite* types are

- lists of native types (except binary objects)
- JavaScript Object Notation (JSON)<sup>1</sup>

*Custom* types include any Erlang term not covered by the previous types. Special care needs to be taken using custom types as they may not be accessible through every API or may be misinterpreted by an API. The use of them is discouraged.

Table 4.2 shows the mapping of supported types to the language-specific types of each API.

<sup>1</sup>see <http://json.org/>



### 4.1.2. JSON API

Scalaris supports a JSON API for transactions. To minimize the necessary round trips between a client and Scalaris, it uses request lists, which contain all requests that can be done in parallel. The request list is then send to a Scalaris node with a POST message. The result contains a list of the results of the requests and - in case of a transaction - a TransLog. To add further requests to the transaction, the TransLog and another list of requests may be send to Scalaris. This process may be repeated as often as necessary. To finish the transaction, the request list can contain a 'commit' request as the last element, which triggers the validation phase of the transaction processing. Request lists are also supported for single read/write operations, i.e. every single operation is committed on its own.

The JSON-API can be accessed via the Scalaris-Web-Server running on port 8000 by default and the page `jsonrpc.yaws` (For example at: <http://localhost:8000/jsonrpc.yaws>). Requests are issued by sending a JSON object with header "Content-type"="application/json" to this URL. The result will then be returned as a JSON object with the same content type. The following table shows how both objects look like:

#### Request

```
{
  "jsonrpc": "2.0",
  "method"  : "<method>",
  "params"  : [<params>],
  "id"      : <number>
}
```

#### Result

```
{
  "result" : <result_object>,
  "id"     : <number>
}
```

The id in the request can be an arbitrary number which identifies the request and is returned in the result. The following operations (shown as `<method>(<params>)`) are currently supported (the given result is the `<result_object>` mentioned above):

- `nop(Value)` - no operation, result:

```
"ok"
```

single read/write:

- `req_list_commit_each(<req_list_rw>)` - commit each request in the list, result:

```
{["status": "ok"] or {"status": "ok", "value": <json_value>} or
 {"status": "fail", "reason": "timeout" or "abort" or "not_found"]}]}
```

- `read(<key>)` - read the value at key, result:

```
{"status": "ok", "value", <json_value>} or
{"status": "fail", "reason": "timeout" or "not_found"}
```

- `write(<key>, <json_value>)` - write value (inside `json_value`) to key, result:

```
{"status": "ok"} or
{"status": "fail", "reason": "timeout" or "abort"}
```

- `test_and_set(<key>, OldValue, NewValue)` - atomic test-and-set (write `NewValue` to key if the current value is `OldValue` - both values are `<json_value>`), result:

```
{ "status": "ok" } or
{ "status": "fail", "reason": "timeout" or "abort" or "not_found" } or
{ "status": "fail", "reason": "key_changed", "value": <json_value> }
```

transactions:

- req\_list(<req\_list>) - process a list of requests, result:

```
{ "tlog": <tlog>,
  "results": [ { "status": "ok" } or { "status": "ok", "value": <json_value> } or
               { "status": "fail", "reason": "timeout" or "abort" or "not_found" } ] }
```

- req\_list(<tlog>, <req\_list>) - process a list of requests with a previous translog, result:

```
{ "tlog": <tlog>,
  "results": [ { "status": "ok" } or { "status": "ok", "value": <json_value> } or
               { "status": "fail", "reason": "timeout" or "abort" or "not_found" } ] }
```

replication layer functions:

- delete(<key>) - delete the value at key, default timeout 2s, result:

```
{ "ok": <number>, "results": [ "ok" or "locks_set" or "undef" ] } or
{ "failure": "timeout", "ok": <number>, "results": [ "ok" or "locks_set" or "undef" ] }
```

- delete(<key>, Timeout) - delete the value at key with a timeout of Timeout Milliseconds, result:

```
{ "ok": <number>, "results": [ "ok" or "locks_set" or "undef" ] } or
{ "failure": "timeout", "ok": <number>, "results": [ "ok" or "locks_set" or "undef" ] }
```

raw DHT functions:

- range\_read(From, To) - read a range of (raw) keys, result:

```
{ "status": "ok" or "timeout",
  "value": [ { "key": <key>, "value": <json_value>, "version": <version> } ] }
```

publish/subscribe:

- publish(Topic, Content) - publish Content to Topic (<key>), result:

```
{ "status": "ok" }
```

- subscribe(Topic, URL) - subscribe URL to Topic (<key>), result:

```
{ "status": "ok" } or
{ "status": "fail", "reason": "timeout" or "abort" }
```

- unsubscribe(Topic, URL) - unsubscribe URL from Topic (<key>), result:

```
{ "status": "ok" } or
{ "status": "fail", "reason": "timeout" or "abort" or "not_found" }
```

- get\_subscribers(Topic) - get subscribers of Topic (<key>), result:

```
[ <urls> ]
```

Note:

```
<json_value> = {"type": "as_is" or "as_bin", "value": <value>}
<req_list_rw> = [{"read", <key>} | {"write", {<key>: <json_value>}}]
<req_list> = [{"read", <key>} | {"write", {<key>: <json_value>}} | {"commit", _}]
```

The <value> inside <json\_value> is either a base64-encoded string representing a binary object (type = "as\_bin") or the value itself (type = "as\_is").

## JSON-Example

The following example illustrates the message flow:

### Client

Make a transaction, that sets two keys →

```
{ "jsonrpc": "2.0",
  "method": "req_list",
  "params": [
    [ { "write": { "keyA": { "type": "as_is", "value": "valueA" } } },
      { "write": { "keyB": { "type": "as_is", "value": "valueB" } } },
      { "commit": "" } ]
  ],
  "id": 0
}
```

### Scalaris node

←

Scalaris sends results back

```
{ "error": null,
  "result": {
    "results": [ { "status": "ok" }, { "status": "ok" }, { "status": "ok" } ],
    "tlog": <TLOG> // this is the translog for further operations!
  },
  "id": 0
}
```

In a second transaction: Read the two keys →

```
{ "jsonrpc": "2.0",
  "method": "req_list",
  "params": [
    [ { "read": "keyA" },
      { "read": "keyB" } ]
  ],
  "id": 0
}
```

←

Scalaris sends results back

```
{ "error": null,
  "result": {
    "results": [
      { "status": "ok", "value": { "type": "as_is", "value": "valueA" } },
      { "status": "ok", "value": { "type": "as_is", "value": "valueB" } }
    ],
    "tlog": <TLOG>
  },
  "id": 0
}
```

Calculate something with the read values →  
 and make further requests, here a write  
 and the commit for the whole transaction. Also include the latest translog we

```
{ "jsonrpc": "2.0",
  "method": "req_list",
  "params": [
    <TLOG>,
    [ { "write": { "keyA": { "type": "as_is", "value": "valueA2" } } },
      { "commit": "" } ]
  ],
  "id": 0
}
```

←

Scalaris sends results back

```
{ "error": null,
  "result": {
    "results": [ { "status": "ok" }, { "status": "ok" } ],
    "tlog": <TLOG>
  },
  "id": 0
}
```

Examples of how to use the JSON API are the Python and Ruby API which use JSON to communicate with Scalaris.

### 4.1.3. Java API

The `scalaris.jar` provides a Java command line client as well as a library for Java programs to access Scalaris. The library provides several classes:

- `TransactionSingleOp` provides methods for reading and writing values.
- `Transaction` provides methods for reading and writing values in transactions.
- `PubSub` provides methods for a simple topic-based pub/sub implementation on top of Scalaris.
- `ReplicatedDHT` provides low-level methods for accessing the replicated DHT of Scalaris.

For details regarding the API we refer the reader to the Javadoc:

```
%> cd java-api
%> ant doc
%> firefox doc/index.html
```

## 4.2. Command Line Interfaces

### 4.2.1. Java command line interface

As mentioned above, the `scalaris.jar` file contains a small command line interface client. For convenience, we provide a wrapper script called `scalaris` which sets up the Java environment:

```
%> ./java-api/scalaris --help
./java-api/scalaris [script options] [options]
```

```

Script Options:
--help, -h          print this message and scalaris help
--noconfig          suppress sourcing of config files in $HOME/.scalaris/
                    and ${prefix}/etc/scalaris/
--execdebug        print scalaris exec line generated by this
                    launch script
--noerl            do not ask erlang for its (local) host name

usage: scalaris [Options]
-b,--minibench <runs> <benchmarks>    run selected mini benchmark(s)
                                         [1|...|9|all] (default: all
                                         benchmarks, 100 test runs)
-d,--delete <key> [<timeout>]          delete an item (default timeout:
                                         2000ms)
                                         WARNING: This function can lead to
                                         inconsistent data (e.g. deleted
                                         items can re-appear). Also when
                                         re-creating an item the version
                                         before the delete can re-appear.
-g,--getsubscribers <topic>            get subscribers of a topic
-h,--help                               print this message
-lh,--localhost                       gets the local host's name as known
                                         to Java (for debugging purposes)
-p,--publish <topic> <message>        publish a new message for the given
                                         topic
-r,--read <key>                        read an item
-s,--subscribe <topic> <url>           subscribe to a topic
-u,--unsubscribe <topic> <url>         unsubscribe from a topic
-v,--verbose                           print verbose information, e.g. the
                                         properties read
-w,--write <key> <value>              write an item

```

read, write and delete can be used to read, write and delete from/to the overlay, respectively. getsubscribers, publish, and subscribe are the PubSub functions. The others provide debugging and testing functionality.

```

%> ./java-api/scalaris -write foo bar
write(foo, bar)
%> ./java-api/scalaris -read foo
read(foo) == bar

```

Per default, the scalaris script tries to connect to a management server at localhost. You can change the node it connects to (and further connection properties) by adapting the values defined in java-api/scalaris.properties.

#### 4.2.2. Python command line interface

```

%> ./python-api/scalaris_client.py --help
usage: ./python-api/scalaris_client.py [Options]
-r,--read <key>                        read an item
-w,--write <key> <value>               write an item
-d,--delete <key> [<timeout>]          delete an item (default timeout:
                                         2000ms)
                                         WARNING: This function can lead to
                                         inconsistent data (e.g. deleted
                                         items can re-appear). Also when
                                         re-creating an item the version
                                         before the delete can re-appear.
-p,--publish <topic> <message>        publish a new message for the given
                                         topic
-s,--subscribe <topic> <url>           subscribe to a topic
-g,--getsubscribers <topic>            get subscribers of a topic
-u,--unsubscribe <topic> <url>         unsubscribe from a topic
-h,--help                               print this message
-b,--minibench <runs> <benchmarks>    run selected mini benchmark(s)

```

```
[1|...|9|all] (default: all  
benchmarks, 100 test runs)
```

### 4.2.3. Ruby command line interface

```
%> ./ruby-api/scalaris_client.rb --help  
Usage: scalaris_client [options]  
  -r, --read KEY           read key KEY  
  -w, --write KEY,VALUE    write key KEY to VALUE  
  -h, --help               Show this message
```

## 5. Testing the system

*Description is based on SVN revision r1618.*

### 5.1. Erlang unit tests

There are some unit tests in the `test` directory which test `Scalaris` itself (the Erlang code). You can call them by running `make test` in the main directory. The results are stored in a local `index.html` file.

The tests are implemented with the `common-test` package from the Erlang system. For running the tests we rely on `run_test`, which is part of the `common-test` package, but (on `erlang < R14`) is not installed by default. `configure` will check whether `run_test` is available. If it is not installed, it will show a warning and a short description of how to install the missing file.

Note: for the unit tests, we are setting up and shutting down several overlay networks. During the shut down phase, the runtime environment will print extensive error messages. These error messages do not indicate that tests failed! Running the complete test suite takes about 10-20 minutes, depending on your machine.

If the test suite is interrupted before finishing, the results may not have been linked into the `index.html` file. They are however stored in the `ct_run.ct@...` directory.

### 5.2. Java unit tests

The Java unit tests can be run by executing `make java-test` in the main directory. This will start a `Scalaris` node with the default ports and test all Java functions part of the Java API. A typical run will look like the following:

```
%> make java-test
[...]
tools.test:
[junit] Running de.zib.tools.PropertyLoaderTest
[junit] Testsuite: de.zib.tools.PropertyLoaderTest
[junit] Tests run: 3, Failures: 0, Errors: 0, Time elapsed: 0.113 sec
[junit] Tests run: 3, Failures: 0, Errors: 0, Time elapsed: 0.113 sec
[junit]
[junit] ----- Standard Output -----
[junit] Working Directory = <scalarisdir>/java-api/classes
[junit] -----
[...]
scalaris.test:
[junit] Running de.zib.scalaris.ConnectionTest
[junit] Testsuite: de.zib.scalaris.ConnectionTest
[junit] Tests run: 7, Failures: 0, Errors: 0, Time elapsed: 0.366 sec
[junit] Tests run: 7, Failures: 0, Errors: 0, Time elapsed: 0.366 sec
[junit]
[junit] Running de.zib.scalaris.DefaultConnectionPolicyTest
[junit] Testsuite: de.zib.scalaris.DefaultConnectionPolicyTest
[junit] Tests run: 12, Failures: 0, Errors: 0, Time elapsed: 0.314 sec
```

```

[junit] Tests run: 12, Failures: 0, Errors: 0, Time elapsed: 0.314 sec
[junit]
[junit] Running de.zib.scalarish.PeerNodeTest
[junit] Testsuite: de.zib.scalarish.PeerNodeTest
[junit] Tests run: 5, Failures: 0, Errors: 0, Time elapsed: 0.077 sec
[junit] Tests run: 5, Failures: 0, Errors: 0, Time elapsed: 0.077 sec
[junit]
[junit] Running de.zib.scalarish.PubSubTest
[junit] Testsuite: de.zib.scalarish.PubSubTest
[junit] Tests run: 33, Failures: 0, Errors: 0, Time elapsed: 4.105 sec
[junit] Tests run: 33, Failures: 0, Errors: 0, Time elapsed: 4.105 sec
[junit]
[junit] ----- Standard Error -----
[junit] 2011-03-25 15:07:04.412: INFO::jetty-7.3.0.v20110203
[junit] 2011-03-25 15:07:04.558: INFO::Started SelectChannelConnector@127.0.0.1:59235
[junit] 2011-03-25 15:07:05.632: INFO::jetty-7.3.0.v20110203
[junit] 2011-03-25 15:07:05.635: INFO::Started SelectChannelConnector@127.0.0.1:41335
[junit] 2011-03-25 15:07:05.635: INFO::jetty-7.3.0.v20110203
[junit] 2011-03-25 15:07:05.643: INFO::Started SelectChannelConnector@127.0.0.1:38552
[junit] 2011-03-25 15:07:05.643: INFO::jetty-7.3.0.v20110203
[junit] 2011-03-25 15:07:05.646: INFO::Started SelectChannelConnector@127.0.0.1:34704
[junit] 2011-03-25 15:07:06.864: INFO::jetty-7.3.0.v20110203
[junit] 2011-03-25 15:07:06.864: INFO::Started SelectChannelConnector@127.0.0.1:57898
[junit] 2011-03-25 15:07:06.864: INFO::jetty-7.3.0.v20110203
[junit] 2011-03-25 15:07:06.865: INFO::Started SelectChannelConnector@127.0.0.1:47949
[junit] 2011-03-25 15:07:06.865: INFO::jetty-7.3.0.v20110203
[junit] 2011-03-25 15:07:06.866: INFO::Started SelectChannelConnector@127.0.0.1:53886
[junit] 2011-03-25 15:07:07.090: INFO::jetty-7.3.0.v20110203
[junit] 2011-03-25 15:07:07.093: INFO::Started SelectChannelConnector@127.0.0.1:33141
[junit] 2011-03-25 15:07:07.094: INFO::jetty-7.3.0.v20110203
[junit] 2011-03-25 15:07:07.096: INFO::Started SelectChannelConnector@127.0.0.1:39119
[junit] 2011-03-25 15:07:07.096: INFO::jetty-7.3.0.v20110203
[junit] 2011-03-25 15:07:07.097: INFO::Started SelectChannelConnector@127.0.0.1:41603
[junit] -----
[junit] Running de.zib.scalarish.ReplicatedDHTTest
[junit] Testsuite: de.zib.scalarish.ReplicatedDHTTest
[junit] Tests run: 6, Failures: 0, Errors: 0, Time elapsed: 0.732 sec
[junit] Tests run: 6, Failures: 0, Errors: 0, Time elapsed: 0.732 sec
[junit]
[junit] Running de.zib.scalarish.TransactionSingleOpTest
[junit] Testsuite: de.zib.scalarish.TransactionSingleOpTest
[junit] Tests run: 28, Failures: 0, Errors: 0, Time elapsed: 0.632 sec
[junit] Tests run: 28, Failures: 0, Errors: 0, Time elapsed: 0.632 sec
[junit]
[junit] Running de.zib.scalarish.TransactionTest
[junit] Testsuite: de.zib.scalarish.TransactionTest
[junit] Tests run: 18, Failures: 0, Errors: 0, Time elapsed: 0.782 sec
[junit] Tests run: 18, Failures: 0, Errors: 0, Time elapsed: 0.782 sec
[junit]

test:

BUILD SUCCESSFUL
Total time: 10 seconds
'jtest_boot@csr-pc9.zib.de'

```

### 5.3. Python unit tests

The Python unit tests can be run by executing `make python-test` in the main directory. This will start a Scalaris node with the default ports and test all Python functions part of the Python API. A typical run will look like the following:

```

%> make python-test
[...]
testDoubleClose (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testRead_NotConnected (TransactionSingleOpTest.TestTransactionSingleOp) ... ok

```



```

testRead_NotFound (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testTestAndSetList1 (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testTestAndSetList2 (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testTestAndSetList_NotConnected (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testTestAndSetList_NotFound (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testTestAndSetString1 (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testTestAndSetString2 (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testTestAndSetString_NotConnected (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testTestAndSetString_NotFound (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testTransactionSingleOp1 (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testTransactionSingleOp2 (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testWriteList1 (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testWriteList2 (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testWriteList_NotConnected (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testWriteString1 (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testWriteString2 (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testWriteString_NotConnected (TransactionSingleOpTest.TestTransactionSingleOp) ... ok
testAbort_Empty (TransactionTest.TestTransaction) ... ok
testAbort_NotConnected (TransactionTest.TestTransaction) ... ok
testCommit_Empty (TransactionTest.TestTransaction) ... ok
testCommit_NotConnected (TransactionTest.TestTransaction) ... ok
testDoubleClose (TransactionTest.TestTransaction) ... ok
testRead_NotConnected (TransactionTest.TestTransaction) ... ok
testRead_NotFound (TransactionTest.TestTransaction) ... ok
testTransaction1 (TransactionTest.TestTransaction) ... ok
testTransaction3 (TransactionTest.TestTransaction) ... ok
testWriteList1 (TransactionTest.TestTransaction) ... ok
testWriteString (TransactionTest.TestTransaction) ... ok
testWriteString_NotConnected (TransactionTest.TestTransaction) ... ok
testWriteString_NotFound (TransactionTest.TestTransaction) ... ok
testDelete1 (ReplicatedDHTTest.TestReplicatedDHT) ... ok
testDelete2 (ReplicatedDHTTest.TestReplicatedDHT) ... ok
testDelete_notExistingKey (ReplicatedDHTTest.TestReplicatedDHT) ... ok
testDoubleClose (ReplicatedDHTTest.TestReplicatedDHT) ... ok
testReplicatedDHT1 (ReplicatedDHTTest.TestReplicatedDHT) ... ok
testReplicatedDHT2 (ReplicatedDHTTest.TestReplicatedDHT) ... ok
testDoubleClose (PubSubTest.TestPubSub) ... ok
testGetSubscribersOtp_NotConnected (PubSubTest.TestPubSub) ... ok
testGetSubscribers_NotExistingTopic (PubSubTest.TestPubSub) ... ok
testPubSub1 (PubSubTest.TestPubSub) ... ok
testPubSub2 (PubSubTest.TestPubSub) ... ok
testPublish1 (PubSubTest.TestPubSub) ... ok
testPublish2 (PubSubTest.TestPubSub) ... ok
testPublish_NotConnected (PubSubTest.TestPubSub) ... ok
testSubscribe1 (PubSubTest.TestPubSub) ... ok
testSubscribe2 (PubSubTest.TestPubSub) ... ok
testSubscribe_NotConnected (PubSubTest.TestPubSub) ... ok
testSubscription1 (PubSubTest.TestPubSub) ... ok
testSubscription2 (PubSubTest.TestPubSub) ... ok
testSubscription3 (PubSubTest.TestPubSub) ... ok
testSubscription4 (PubSubTest.TestPubSub) ... ok
testUnsubscribe1 (PubSubTest.TestPubSub) ... ok
testUnsubscribe2 (PubSubTest.TestPubSub) ... ok
testUnsubscribe_NotConnected (PubSubTest.TestPubSub) ... ok
testUnsubscribe_NotExistingTopic (PubSubTest.TestPubSub) ... ok
testUnsubscribe_NotExistingUrl (PubSubTest.TestPubSub) ... ok

```

```

-----
Ran 58 tests in 12.317s

```

```
OK
```

```
'jtest_boot@csr-pc9.zib.de'
```

## 5.4. Interoperability Tests

In order to check whether the common types described in Section 4.1 on page 15 are fully supported by the APIs and yield to the appropriate types in another API, we implemented some interoperability tests. They can be run by executing `make interop-test` in the main directory. This will start a Scalaris node with the default ports, write test data using both the Java and the Python APIs and let each API read the data it wrote itself as well as the data the other API read. On success it will print

```
%> make interop-test  
[...]  
all tests successful
```

## 6. Troubleshooting

*Description is based on SVN revision r1618.*

### 6.1. Network

Scalaris uses a couple of TCP ports for communication. It does not use UDP at the moment.

	HTTP Server	Inter-node communication
default (see bin/scalaris.cfg)	8000	14195–14198
first node (bin/firstnode.sh)	8000	14195
joining node 1 (bin/joining_node.sh)	8001	14196
other joining nodes (bin/joining_node.sh <ID>)	8000 + <ID>	14195 + <ID>
standalone mgmt server (bin/mgmt-server.sh)	7999	14194

Please make sure that at least 14195 and 14196 are not blocked by firewalls in order to be able to start at least one first and one joining node on each machine..

### 6.2. Miscellaneous

For up-to-date information about frequently asked questions and troubleshooting, please refer to our FAQs at <https://code.google.com/p/scalaris/wiki/FAQ> and our mailing list at <http://groups.google.com/group/scalaris>.

## Part II.

# Developers Guide

## 7. General Hints

### 7.1. Coding Guidelines

- Keep the code short
- Use `gen_component` to implement additional processes
- Don't use `receive` by yourself (Exception: to implement single threaded user API calls (`cs_api`, `yaws_calls`, etc))
- Don't use `erlang:now/0`, `erlang:send_after/3`, `receive after` etc. in performance critical code, consider using `msg_delay` instead.
- Don't use `timer:tc/3` as it catches exceptions. Use `util:tc/3` instead.

### 7.2. Testing Your Modifications and Extensions

- Run the testsuites using `make test`
- Run the java api test using `make java-test` (Scalaris output will be printed if a test fails; if you want to see it during the tests, start a `bin/firstnode.sh` and run the tests by `cd java; ant test`)
- Run the Ruby client by starting Scalaris and running `cd contrib; ./jsonrpc.rb`

### 7.3. Help with Digging into the System

- use `ets:i/0,1` to get details on the local state of some processes
- consider changing `pdb.erl` to use `ets` instead of `erlang:put/get`
- Have a look at `strace -f -p PID` of beam process
- Get message statistics via the Web-interface
- enable/disable tracing for certain modules
- Use `etop` and look at the total memory size and atoms generated
- send processes `sleep` or `kill` messages to test certain behaviour (see `gen_component.erl`)
- use `mgmt_server:number_of_nodes(). flush().`
- use `admin_checkring(). flush().`

## 8. System Infrastructure

### 8.1. Groups of Processes

- What is it? How to distinguish from Erlangs internal named processes?
- Joining a process group
- Why do we do this... (managing several independent nodes inside a single Erlang VM for testing)

### 8.2. The Communication Layer `comm`

- in general
- format of messages (tuples)
- use messages with cookies (server and client side)
- What is a message tag?

### 8.3. The `gen_component`

*Description is based on SVN revision r1620.*

The generic component model implemented by `gen_component` allows to add some common functionality to all the components that build up the Scalaris system. It supports:

**event-handlers:** message handling with a similar syntax as used in [3].

**FIFO order of messages:** components cannot be inadvertently locked as we do not use selective receive statements in the code.

**sleep and halt:** for testing components can sleep or be halted.

**debugging, breakpoints, stepwise execution:** to debug components execution can be steered via breakpoints, step-wise execution and continuation based on arriving events and user defined component state conditions.

**basic profiling,**

**state dependent message handlers:** depending on its state, different message handlers can be used and switched during runtime. Thereby a kind of state-machine based message handling is supported.

**prepared for `pid_groups`:** allows to send events to named processes inside the same group as the actual component itself (`send_to_group_member`) when just holding a reference to any group member, and

**unit-testing of event-handlers:** as message handling is separated from the main loop of the component, the handling of individual messages and thereby performed state manipulation can easily be tested in unit-tests by directly calling message handlers.

In Scalaris all Erlang processes should be implemented as `gen_component`. The only exception are functions interfacing to the client, where a transition from asynchronous to synchronous request handling is necessary and that are executed in the context of a client's process or a process that behaves as a proxy for a client (`cs_api`).

### 8.3.1. A basic `gen_component` including a message handler

To implement a `gen_component`, the component has to provide the `gen_component` behaviour:

File `gen_component.erl`:

```

83 -spec behaviour_info(atom()) -> [{atom(), arity()}] | undefined.
84 behaviour_info(callbacks) ->
85   [
86     {init, 1}      % initialize component
87     % note: can use arbitrary on-handler, but by default on/2 is used:
88     %% {on, 2}      % handle a single message
89     %% {on(Msg, State) -> NewState | unknown_event | kill
90   ];

```

This is illustrated by the following example:

File `msg_delay.erl`:

```

70 %% initialize: return initial state.
71 -spec init([]) -> state().
72 init([]) ->
73   MyGroup = pid_groups:my_groupname(),
74   ?TRACE("msg_delay:init for pid group ~p~n", [MyGroup]),
75   TimeTableName = list_to_atom(MyGroup ++ "_msg_delay"),
76   %% use random table name provided by ets to *not* generate an atom
77   %% TableName = pdb:new(?MODULE, [set, private]),
78   TimeTable = pdb:new(TimeTableName, [set, protected, named_table]),
79   comm:send_local(self(), {msg_delay_periodic}),
80   _State = {TimeTable, _Round = 0}.
81
82 -spec on(message(), state()) -> state().
83 on({msg_delay_req, Seconds, Dest, Msg} = _FullMsg,
84   {TimeTable, Counter} = State) ->
85   ?TRACE("msg_delay:on(~.0p, ~.0p)~n", [_FullMsg, State]),
86   Future = trunc(Counter + Seconds),
87   case pdb:get(Future, TimeTable) of
88     undefined ->
89       pdb:set({Future, [{Dest, Msg}]}, TimeTable);
90     {_, MsgQueue} ->
91       pdb:set({Future, [{Dest, Msg} | MsgQueue]}, TimeTable)
92   end,
93   State;
94
95 %% periodic trigger
96 on({msg_delay_periodic} = Trigger, {TimeTable, Counter} = _State) ->
97   ?TRACE("msg_delay:on(~.0p, ~.0p)~n", [Trigger, State]),
98   case pdb:get(Counter, TimeTable) of
99     undefined -> ok;
100    {_, MsgQueue} ->
101      _ = [ comm:send_local(Dest, Msg) || {Dest, Msg} <- MsgQueue ],
102      pdb:delete(Counter, TimeTable)
103   end,
104   comm:send_local_after(1000, self(), Trigger),
105   {TimeTable, Counter + 1};
106
107 on({web_debug_info, Requestor}, {TimeTable, Counter} = State) ->
108   KeyValueCollection =
109     [{"queued messages (in 0-10s, messages):", ""}] |
110     [begin
111       Future = trunc(Counter + Seconds),

```

```

112         Queue = case pdb:get(Future, TimeTable) of
113             undefined -> none;
114             {_, Q}    -> Q
115         end,
116         {lists:flatten(io_lib:format("~p", [Seconds])),
117          lists:flatten(io_lib:format("~p", [Queue]))}
118     end || Seconds <- lists:seq(0, 10)]],
119     comm:send_local(Requestor, {web_debug_info_reply, KeyValueType}),
120     State.

```

`your_gen_component:init/1` is called during start-up of a `gen_component` and should return the initial state to be used for this `gen_component`. Later, the current state of the component can be retrieved using `gen_component:get_state/1`.

To react on messages / events, a message handler is used. The default message handler is called `your_gen_component:on/2`. This can be changed by calling `gen_component:change_handler/2` (see Section 8.3.6). When an event / message for the component arrives, this handler is called with the event itself and the current state of the component. In the handler, the state of the component may be adjusted depending upon the event. The handler itself may trigger new events / messages for itself or other components and has finally to return the updated state of the component or the atoms `unknown_event` or `kill`. It must neither call `receive` nor `timer:sleep/1` nor `erlang:exit/1`.

### 8.3.2. How to start a `gen_component`?

A `gen_component` can be started using one of:

```
gen_component:start(Module, Args, GenCOptions = [])
```

```
gen_component:start_link(Module, Args, GenCOptions = [])
```

Module: the name of the module your component is implemented in

Args: List of parameters passed to `Module:init/1` for initialization

GenCOptions: optional parameter. List of options for `gen_component`

`{pid_groups_join_as, ProcessGroup, ProcessName}`: registers the new process with the given process group (also called `instanceid`) and name using `pid_groups`.

`{erlang_register, ProcessName}`: registers the process as a named Erlang process.

`wait_for_init`: wait for `Module:init/1` to return before returning to the caller.

These functions are compatible to the Erlang/OTP supervisors. They spawn a new process for the component which itself calls `Module:init/1` with the given `Args` to initialize the component. `Module:init/1` should return the initial state for your component. For each message sent to this component, the default message handler `Module:on(Message, State)` will be called, which should react on the message and return the updated state of your component.

`gen_component:start()` and `gen_component:start_link()` return the pid of the spawned process as `{ok, Pid}`.

### 8.3.3. When does a `gen_component` terminate?

A `gen_component` can be stopped using:

`gen_component:kill(Pid)` or by returning `kill` from the current message handler.



### 8.3.4. What happens when unexpected events / messages arrive?

Your message handler (default is `your_gen_component:on/2`) should return `unknown_event` in the final clause (`your_gen_component:on(_,_)`). `gen_component` then will nicely report on the unhandled message, the component's name, its state and currently active message handler, as shown in the following example:

```
# bin/boot.sh
[...]
(boot@localhost)10> pid_groups ! {no_message}.
{no_message}
[error] unknown message: {no_message} in Module: pid_groups and
handler on in State null
(boot@localhost)11>
```

The `pid_groups` (see Section 8.1) is a `gen_component` which registers itself as named Erlang process with the `gen_component` option `erlang_register` and therefore can be addressed by its name in the Erlang shell. We send it a `{no_message}` and `gen_component` reports on the unhandled message. The `pid_groups` module itself continues to run and waits for further messages.

### 8.3.5. What if my message handler generates an exception or crashes the process?

`gen_component` catches exceptions generated by message handlers and reports them with a stack trace, the message, that generated the exception, and the current state of the component.

If a message handler terminates the process via `erlang:exit/1`, this is out of the responsibility scope of `gen_component`. As usual in Erlang, all linked processes will be informed. If for example `gen_component:start_link/2` or `/3` was used for starting the `gen_component`, the spawning process will be informed, which may be an Erlang supervisor process taking further actions.

### 8.3.6. Changing message handlers and implementing state dependent message responsiveness as a state-machine

Sometimes it is beneficial to handle messages depending on the state of a component. One possibility to express this is implementing different clauses depending on the state variable, another is introducing case clauses inside message handlers to distinguish between current states. Both approaches may become tedious, error prone, and may result in confusing source code.

Sometimes the use of several different message handlers for different states of the component leads to clearer arranged code, especially if the set of handled messages changes from state to state. For example, if we have a component with an initialization phase and a production phase afterwards, we can handle in the first message handler messages relevant during the initialization phase and simply queue all other requests for later processing using a common default clause.

When initialization is done, we handle the queued user requests and switch to the message handler for the production phase. The message handler for the initialization phase does not need to know about messages occurring during production phase and the message handler for the production phase does not need to care about messages used during initialization. Both handlers can be made independent and may be extended later on without any adjustments to the other.

One can also use this scheme to implement complex state-machines by changing the message handler from state to state.

To switch the message handler `gen_component:change_handler(State, new_handler)` is called as

the last operation after a message in the active message handler was handled, so that the return value of `gen_component:change_handler/2` is propagated to `gen_component`. The new handler is given as an atom, which is the name of the 2-ary function in your component module to be called.

### Starting with non-default message handler.

It is also possible to change the message handler right from the start in your `your_gen_component:init/1` to avoid the default message handler `your_gen_component:on/2`. Just create your initial state as usual and call `gen_component:change_handler(State, my_handler)` as the final call in your `your_gen_component:init/1`. We prepared `gen_component:change_handler/2` to return `State` itself, so this will work properly.

### 8.3.7. Handling several messages atomically

The message handler is called for each message separately. Such a single call is atomic, i.e. the component does not perform any other action until the called message handler finishes. Sometimes, it is necessary to execute two or more calls to the message handler atomically (without other interleaving messages). For example if a message A contains another message B as payload, it may be necessary to handle A and B directly one after the other without interference of other messages. So, after handling A you want to call your message handler with B.

In most cases, you could just do so by calculating the new state as result of handling message A first and then calling the message handler with message B and the new state by yourself.

It is safer to use `gen_component:post_op(2)` in such cases: When *B* contains a special message, which is usually handled by the `gen_component` module itself (like `send_to_group_member`, `kill`, `sleep`), the direct call to the message handler would not achieve the expected result. By calling `gen_component:post_op(NewState, B)` to return the new state after handling message A, message B will be handled directly after the current message A.

### 8.3.8. Halting and pausing a gen\_component

Using `gen_component:kill(Pid)` and `gen_component:sleep(Pid, Time)` components can be terminated or paused.

### 8.3.9. Integration with pid\_groups: Redirecting messages to other gen\_components

Each `gen_component` by itself is prepared to support `comm:send_to_group_member/3` which forwards messages inside a group of processes registered via `pid_groups` (see Section 8.1) by their name. So, if you hold a `Pid` of one member of a process group, you can send messages to other members of this group, if you know their registered Erlang name. You do not necessarily have to know their individual `Pid`.

*In consequence, no `gen_component` can individually handle messages of the form `{send_to_group_member, _, _}` as such messages are consumed by `gen_component` itself.*

### 8.3.10. Replying to ping messages

Each `gen_component` replies automatically to `{ping, Pid}` requests with a `{pong}` send to the given `Pid`. Such messages are generated, for example, by `vivaldi_latency` which is used by our `vivaldi` module.

*In consequence, no `gen_component` can individually handle messages of the form: `{ping, _}` as such messages are consumed by `gen_component` itself.*

### 8.3.11. The debugging interface of `gen_component`: Breakpoints and step-wise execution

We equipped `gen_component` with a debugging interface, which especially is beneficial, when testing the interplay between several `gen_components`. It supports breakpoints (bp) which can pause the `gen_component` depending on the arriving messages or depending on user defined conditions. If a breakpoint is reached, the execution can be continued step-wise (message by message) or until the next breakpoint is reached.

We use it in our unit tests to steer protocol interleavings and to perform tests using random protocol interleavings between several processes (see `paxos_SUITE`). It allows also to reproduce given protocol interleavings for better testing.

#### Managing breakpoints.

Breakpoints are managed by the following functions:

`gen_component:bp_set(Pid, MsgTag, BPName)`: For the component running under `Pid` a breakpoint `BPName` is set. It is reached, when a message with a message tag `MsgTag` is next to be handled by the component (See `comm:get_msg_tag/1` and Section 8.2 for more information on message tags). The `BPName` is used as a reference for this breakpoint, for example to delete it later.

`gen_component:bp_set_cond(Pid, Cond, BPName)`: The same as `gen_component:bp_set/3` but a user defined condition implemented in `{Module, Function, Params = 2}` = `Cond` is checked by calling `Module:Function(Message, State)` to decide whether a breakpoint is reached or not. `Message` is the next message to be handled by the component and `State` is the current state of the component. `Module:Function/2` should return a boolean.

`gen_component:bp_del(Pid, BPName)`: The breakpoint `BPName` is deleted. If the component is in this breakpoint, it will not be released by this call. This has to be done separately by `gen_component:bp_cont/1`. But the deleted breakpoint will no longer be considered for newly entering a breakpoint.

`gen_component:bp_barrier(Pid)`: Delay all further handling of breakpoint requests until a breakpoint is actually entered.

*Note, that the following call sequence may not catch the breakpoint at all, as during the sleep the component not necessarily consumes a ping message and the set breakpoint 'sample\_bp' may already be deleted before a ping message arrives.*

```
gen_component:bp_set(Pid, ping, sample_bp),
timer:sleep(10),
gen_component:bp_del(Pid, sample_bp),
gen_component:bp_cont(Pid).
```

To overcome this, `gen_component:bp_barrier/1` can be used:

```
gen_component:bp_set(Pid, ping, sample_bp),
gen_component:bp_barrier(Pid),
%% After the bp_barrier request, following breakpoint requests
%% will not be handled before a breakpoint is actually entered.
%% The gen_component itself is still active and handles messages as usual
%% until it enters a breakpoint.
gen_component:bp_del(Pid, sample_bp),
% Delete the breakpoint after it was entered once (ensured by bp_barrier).
% Release the gen_component from the breakpoint and continue.
gen_component:bp_cont(Pid).
```

None of the calls in the sample listing above is blocking. It just schedules all the operations, including the `bp_barrier`, for the `gen_component` and immediately finishes. The actual events of entering and continuing the breakpoint in the `gen_component` happens independently later on, when the next ping message arrives.

## Managing execution.

The execution of a `gen_component` can be managed by the following functions:

`gen_component:bp_step(Pid)`: This is the only blocking breakpoint function. It waits until the `gen_component` is in a breakpoint and has handled a single message. It returns the module, the active message handler, and the handled message as a tuple `{Module, On, Message}`. This function does not actually finish the breakpoint, but just lets a single message pass through. For further messages, no breakpoint condition has to be valid, the original breakpoint is still active. To leave a breakpoint, use `gen_component:bp_cont/1`.

`gen_component:bp_cont(Pid)`: Leaves a breakpoint. `gen_component` runs as usual until the next breakpoint is reached.

If no further breakpoints should be entered after continuation, you should delete the registered breakpoint using `gen_component:bp_del/2` before continuing the execution with `gen_component:bp_cont/1`. To ensure, that the breakpoint is entered at least once, `gen_component:bp_barrier/1` should be used before deleting the breakpoint (see the example above). Otherwise it could happen, that the delete request arrives at your `gen_component` before it was actually triggered. The following continuation request would then unintentional apply to an unrelated breakpoint that may be entered later on.

`gen_component:runnable(Pid)`: Returns whether a `gen_component` has messages to handle and is runnable. If you know, that a `gen_component` is in a breakpoint, you can use this to check, whether a `gen_component:bp_step/1` or `gen_component:bp_cont/1` is applicable to the component.

## Tracing handled messages – getting a message interleaving protocol.

We use the debugging interface of `gen_component` to test protocols with random interleaving. First we start all the components involved, set breakpoints on the initialization messages for a new Paxos consensus and then start a single Paxos instance on all of them. The outcome of the Paxos consensus is a `learner_decide` message. So, in `paxos_SUITE:step_until_decide/3` we look for runnable processes and select randomly one of them to perform a single step until the protocol finishes with a decision.

File paxos\_SUITE.erl:

```

236 -spec prop_rnd_interleave(1..4, 4..16, {pos_integer(), pos_integer(), pos_integer()})
237     -> true.
238 prop_rnd_interleave(NumProposers, NumAcceptors, Seed) ->
239     ct:pal("Called with: paxos_SUITE:prop_rnd_interleave(~p, ~p, ~p).~n",
240         [NumProposers, NumAcceptors, Seed]),
241     Majority = NumAcceptors div 2 + 1,
242     {Proposers, Acceptors, Learners} =
243         make(NumProposers, NumAcceptors, 1, "rnd_interleave"),
244     %% set bp on all processes
245     _ = [ gen_component:bp_set(comm:make_local(X), proposer_initialize, bp)
246         || X <- Proposers ],
247     _ = [ gen_component:bp_set(comm:make_local(X), acceptor_initialize, bp)
248         || X <- Acceptors ],
249     _ = [ gen_component:bp_set(comm:make_local(X), learner_initialize, bp)
250         || X <- Learners ],
251     %% start paxos instances
252     _ = [ proposer:start_paxosid(X, paxidrndinterl, Acceptors,
253         proposal, Majority, NumProposers, Y)
254         || {X,Y} <- lists:zip(Proposers, lists:seq(1, NumProposers)) ],
255     _ = [ acceptor:start_paxosid(X, paxidrndinterl, Learners)
256         || X <- Acceptors ],
257     _ = [ learner:start_paxosid(X, paxidrndinterl, Majority,
258         comm:this(), cpaxidrndinterl)
259         || X <- Learners ],
260     %% randomly step through protocol
261     OldSeed = random:seed(Seed),
262     Steps = step_until_decide(Proposers ++ Acceptors ++ Learners, cpaxidrndinterl, 0),
263     ct:pal("Needed ~p steps~n", [Steps]),
264     _ = case OldSeed of
265         undefined -> ok;
266         _ -> random:seed(OldSeed)
267     end,
268     _ = [ gen_component:kill(comm:make_local(X))
269         || X <- lists:flatten([Proposers, Acceptors, Learners]) ],
270     true.
271
272 step_until_decide(Processes, PaxId, SumSteps) ->
273     %% io:format("Step ~p~n", [SumSteps]),
274     Runnable = [ X || X <- Processes, gen_component:runnable(comm:make_local(X)) ],
275     case Runnable of
276     [] ->
277         ct:pal("No runnable processes of ~p~n", [length(Processes)]),
278         timer:sleep(5), step_until_decide(Processes, PaxId, SumSteps);
279     _ ->
280         Num = random:uniform(length(Runnable)),
281         _ = gen_component:bp_step(comm:make_local(lists:nth(Num, Runnable))),
282         receive
283             {learner_decide, cpaxidrndinterl, _, _Res} = _Any ->
284                 %% io:format("Received ~p~n", [_Any]),
285                 SumSteps
286         after 0 -> step_until_decide(Processes, PaxId, SumSteps + 1)
287         end
288     end.

```

To get a message interleaving protocol, we either can output the results of each `gen_component:bp_step/1` call together with the `Pid` we selected for stepping, or alter the definition of the macro `TRACE_BP_STEPS` in `gen_component`, when we execute all `gen_components` locally in the same Erlang virtual machine.

File `gen_component.erl`:

```

31 %-define(TRACE_BP_STEPS(X,Y), io:format(X,Y)).           %% output on console
32 %-define(TRACE_BP_STEPS(X,Y), ct:pal(X,Y)).             %% output even if called by unittest
33 %-define(TRACE_BP_STEPS(X,Y), io:format(user,X,Y)).     %% clean output even if called by unittest
34 -define(TRACE_BP_STEPS(X,Y), ok).

```

### 8.3.12. Future use and planned extensions for `gen_component`

`gen_component` could be further extended. For example it could support hot-code upgrade or could be used to implement algorithms that have to be run across several components of Scalaris like snapshot algorithms or similar extensions.

## 8.4. The Process' Database (`pdb`)

- How to use it and how to switch from `erlang:put/set` to `ets` and implied limitations.

## 8.5. Failure Detectors (`fd`)

- uses Erlang monitors locally
- is independent of component load
- uses heartbeats between Erlang virtual machines
- uses a single proxy heartbeat server per Erlang virtual machine, which itself uses Erlang monitors to monitor locally
- uses dynamic timeouts to implement an eventually perfect failure detector.

## 8.6. Monitoring Statistics (`monitor`, `rrd`)

*Description is based on SVN revision r2546.*

The `monitor` module offers several methods to gather meaningful statistics using the `rrd()` data type defined in `rrd`.

`rrd()` records work with time slots, i.e. a fixed slot length is given at creation and items which should be inserted will be either put into the current slot, or a new slot will be created. Each data item thus needs a time stamp associated with it. It must not be a real time, but can also be a virtual time stamp.

The `rrd` module thus offers two different APIs: one with transparent time handling, e.g. `rrd:create/3`, `rrd:add_now/2`, and one with manual time handling, e.g. `rrd:create/4`, `rrd:add/3`.

To allow different evaluations of the stored data, the following types of data are supported:

- `gauge`: only stores the newest value of a time slot, e.g. for thermometers,
- `counter`: sums up all values inside a time slot,
- `timing`: records time spans and stores values to easily calculate e.g. the sum, the standard deviation, the number of events, the min and max,
- `timing_with_hist`: similar to `timing` but also records a more detailed (approximated) histogram of the data,
- `event`: records each event (including its time stamp) inside a time slot in a list (this should be rarely used as the amount of data stored may be very big).

The monitor offers functions to conveniently store and retrieve such values. It is also started as a process in each `dht_node` and `basic_services` group as well as inside each `clients_group`. This process ultimately stores the whole `rrd()` structure. There are three paradigms how values can be stored:

1. Values are gathered in the process that is generating the values. Inside this process, the `rrd()` is stored in the erlang dictionary. Whenever a new time slot is started, the values will be reported to the monitor process of the gathering process' group.
2. Values are gathered in the process that is generating the values. Inside this process, the `rrd()` is handled manually. After changing the `rrd()`, a manual check for reporting needs to be issued using `monitor:check_report/4`.
3. Values are immediately send to the monitor process where it undergoes the same procedures until it is finally stored and available to other processes. This is especially useful if the process generating the values does not live long or does not regularly create new data, e.g. the client.

The following example illustrates the first mode, i.e. gathering data in the generating process. It has been taken from the `cyclon` module which uses a counter data type:

```
% initialise the monitor with an empty rrd() using a 60s monitoring interval
monitor:proc_set_value(?MODULE, 'shuffle', rrd:create(60 * 1000000, 3, counter)),
% update the value by adding one
monitor:proc_set_value(?MODULE, 'shuffle', fun(Old) -> rrd:add_now(1, Old) end),
% check regularly whether to report the data to the monitor:
monitor:proc_check_timeslot(?MODULE, 'shuffle')
```

The first two parameters of `monitor:proc_set_value/3` define the name of a monitored value, the module's name and a unique key. The second can be either an `rrd()` or an update fun. The `monitor:proc_check_timeslot/3` function can be used if your module does not regularly create new data. In this case, the monitor process would not have the latest data for others to retrieve. This function forces a check and creates the new time slot if needed (thus reporting the data).

This is how forwarding works (taken from `api_tx`):

```
monitor:client_monitor_set_value(
  ?MODULE, 'req_list',
  fun(Old) ->
    Old2 = case Old of
      % 10s monitoring interval, only keep newest in the client process
      undefined -> rrd:create(10 * 1000000, 1, {timing, ms});
      _ -> Old
    end,
    rrd:add_now(TimeInUs / 1000, Old2)
  end),
```

As in this case there is no safe way of initialising the value, it is more useful to provide an update fun to `monitor:client_monitor_set_value/3`. This function is only useful for the client processes as it reports to the monitor in the `clients_group` (recall that client processes do not belong to any group). All other processes should use `monitor:monitor_set_value/3` with the same semantics.

## 8.7. Writing Unittests

### 8.7.1. Plain unittests

### 8.7.2. Randomized Testing using `tester.erl`



## 9. Basic Structured Overlay

### 9.1. Ring Maintenance

### 9.2. T-Man

### 9.3. Routing Tables

*Description is based on SVN revision r1453.*

Each node of the ring can perform searches in the overlay.

A search is done by a lookup in the overlay, but there are several other demands for communication between peers. Scalaris provides a general interface to route a message to the (other) peer, which is currently responsible for a given key.

File `api_dht_raw.erl`:

```
31 -spec unreliable_lookup(Key::?RT:key(), Msg::comm:message()) -> ok.
32 unreliable_lookup(Key, Msg) ->
33     comm:send_local(pid_groups:find_a(dht_node),
34                     {lookup_aux, Key, 0, Msg}).
35
36 -spec unreliable_get_key(Key::?RT:key()) -> ok.
37 unreliable_get_key(Key) ->
38     unreliable_lookup(Key, {get_key, comm:this(), Key}).
39
40 -spec unreliable_get_key(CollectorPid::comm:mypid(),
41                         ReqId::{rdht_req_id, pos_integer()},
42                         Key::?RT:key()) -> ok.
43 unreliable_get_key(CollectorPid, ReqId, Key) ->
44     unreliable_lookup(Key, {get_key, CollectorPid, ReqId, Key}).
```

The message `Msg` could be a `get_key` which retrieves content from the responsible node or a `get_node` message, which returns a pointer to the node.

All currently supported messages are listed in the file `dht_node.erl`.

The message routing is implemented in `dht_node_lookup.erl`

File `dht_node_lookup.erl`:

```
27 %% @doc Find the node responsible for Key and send him the message Msg.
28 -spec lookup_aux(State::dht_node_state:state(), Key::intervals:key(),
29                 Hops::non_neg_integer(), Msg::comm:message()) -> ok.
30 lookup_aux(State, Key, Hops, Msg) ->
31     Neighbors = dht_node_state:get(State, neighbors),
32     case intervals:in(Key, nodelist:succ_range(Neighbors)) of
33     true -> % found node -> terminate
34         P = node:pidX(nodelist:succ(Neighbors)),
35         comm:send(P, {lookup_fin, Key, Hops + 1, Msg}, [{shepherd, self()}]);
36     ->
37         P = ?RT:next_hop(State, Key),
38         comm:send(P, {lookup_aux, Key, Hops + 1, Msg}, [{shepherd, self()}])
```



```
39     end.
```

Each node is responsible for a certain key interval. The function `intervals:in/2` is used to decide, whether the key is between the current node and its successor. If that is the case, the final step is delivers a `lookup_fin` message to the local node. Otherwise, the message is forwarded to the next nearest known peer (listed in the routing table) determined by `?RT:next_hop/2`.

`rt_beh.erl` is a generic interface for routing tables. It can be compared to interfaces in Java. In Erlang interfaces can be defined using a so called ‘behaviour’. The files `rt_simple` and `rt_chord` implement the behaviour ‘`rt_beh`’.

The macro `?RT` is used to select the current implementation of routing tables. It is defined in `include/scalaris.hrl`.

File `scalaris.hrl`:

```
28 %%The RT macro determines which kind of routingtable is used. Uncomment the
29 %%one that is desired.
30
31 %%Standard Chord routingtable
32 -define(RT, rt_chord).
33 % first valid key:
34 -define(MINUS_INFINITY, 0).
35 -define(MINUS_INFINITY_TYPE, 0).
36 % first invalid key:
37 -define(PLUS_INFINITY, 16#10000000000000000000000000000000).
38 -define(PLUS_INFINITY_TYPE, 16#10000000000000000000000000000000).
39
40 %%Simple routingtable
41 %-define(RT, rt_simple).
```

The functions, that have to be implemented for a routing mechanism are defined in the following file:

File `rt_beh.erl`:

```
32 -spec behaviour_info(atom()) -> [{atom(), arity()}] | undefined.
33 behaviour_info(callbacks) ->
34 [
35     % create a default routing table
36     {empty, 1}, {empty_ext, 1},
37     % mapping: key space -> identifier space
38     {hash_key, 1}, {get_random_node_id, 0},
39     % routing
40     {next_hop, 2},
41     % trigger for new stabilization round
42     {init_stabilize, 2},
43     % adapt RT to changed neighborhood
44     {update, 3},
45     % dead nodes filtering
46     {filter_dead_node, 2},
47     % statistics
48     {to_pid_list, 1}, {get_size, 1},
49     % gets all (replicated) keys for a given (hashed) key
50     % (for symmetric replication)
51     {get_replica_keys, 1},
52     % address space size, range and split key
53     % (may all throw 'throw:not_supported' if unsupported by the RT)
54     {n, 0}, {get_range, 2}, {get_split_key, 3},
55     % for debugging and web interface
56     {dump, 1},
57     % for bulkowner
58     {to_list, 1},
59     % convert from internal representation to version for dht_node
60     {export_rt_to_dht_node, 2},
61     % handle messages specific to a certain routing-table implementation
```

```

62     {handle_custom_message, 2},
63     % common methods
64     {check, 4}, {check, 5},
65     {check_config, 0}
66 ];

```

empty/1 gets a successor and generates an empty routing table for use inside the routing table implementation. The data structure of the routing table is undefined. It can be a list, a tree, a matrix ...

empty\_ext/1 similarly creates an empty external routing table for use by the dht\_node. This process might not need all the information a routing table implementation requires and can thus work with less data.

hash\_key/1 gets a key and maps it into the overlay's identifier space.

get\_random\_node\_id/0 returns a random node id from the overlay's identifier space. This is used for example when a new node joins the system.

next\_hop/2 gets a dht\_node's state (including the external routing table representation) and a key and returns the node, that should be contacted next when searching for the key, i.e. the known node nearest to the id.

init\_stabilize/2 is called periodically to rebuild the routing table. The parameters are the identifier of the node, its successor and the old (internal) routing table state. This method may send messages to the routing\_table process which need to be handled by the handle\_custom\_message/handler since they are implementation-specific.

update/7 is called when the node's ID, predecessor and/or successor changes. It updates the (internal) routing table with the (new) information.

filter\_dead\_node/2 is called by the failure detector and tells the routing table about dead nodes. This function gets the (internal) routing table and a node to remove from it. A new routing table state is returned.

to\_pid\_list/1 get the PIDs of all (internal) routing table entries.

get\_size/1 get the (internal or external) routing table's size.

get\_replica\_keys/1 Returns for a given (hashed) Key the (hashed) keys of its replicas. This used for implementing symmetric replication.

n/0 gets the number of available keys. An implementation may throw `throw:not_supported` if the operation is unsupported by the routing table.

dump/1 dump the (internal) routing table state for debugging, e.g. by using the web interface. Returns a list of `{Index, Node_as_String}` tuples which may just as well be empty.

to\_list/1 convert the (external) representation of the routing table inside a given dht\_node\_state to a sorted list of known nodes from the routing table, i.e. first=succ, second=next known node on the ring, ... This is used by bulk-operations to create a broadcast tree.

export\_rt\_to\_dht\_node/2 convert the internal routing table state to an external state. Gets the internal state and the node's neighborhood for doing so.

handle\_custom\_message/2 handle messages specific to the routing table implementation. rt\_loop will forward unknown messages to this function.

check/5, check/6 check for routing table changes and send an updated (external) routing table to the dht\_node process.

check\_config/0 check that all required configuration parameters exist and satisfy certain restrictions.

### 9.3.1. The routing table process (rt\_loop)

The `rt_loop` module implements the process for all routing tables. It processes messages and calls the appropriate methods in the specific routing table implementations.

File `rt_loop.erl`:

```
40 -opaque(state_active() :: {Neighbors      :: nodelist:neighborhood(),
41                               RTState      :: ?RT:rt(),
42                               TriggerState  :: trigger:state()}).
43 -type(state_inactive() :: {inactive,
44                               MessageQueue::msg_queue:msg_queue(),
45                               TriggerState::trigger:state()}).
46 %% -type(state() :: state_active() | state_inactive()).
```

If initialized, the node's id, its predecessor, successor and the routing table state of the selected implementation (the macro `RT` refers to).

File `rt_loop.erl`:

```
153 on_active({trigger_rt}, {Neighbors, OldRT, TriggerState}) ->
154     % start periodic stabilization
155     % log:log(debug, "[ RT ] stabilize"),
156     NewRT = ?RT:init_stabilize(Neighbors, OldRT),
157     ?RT:check(OldRT, NewRT, Neighbors, true),
158     % trigger next stabilization
159     NewTriggerState = trigger:next(TriggerState),
160     new_state(Neighbors, NewRT, NewTriggerState);
```

Periodically (see `routingtable_trigger` and `pointer_base_stabilization_interval` config parameters) a trigger message is sent to the `rt_loop` process that starts the periodic stabilization implemented by each routing table.

File `rt_loop.erl`:

```
138 % update routing table with changed ID, pred and/or succ
139 on_active({update_rt, OldNeighbors, NewNeighbors}, {_Neighbors, OldRT, TriggerState}) ->
140     case ?RT:update(OldRT, OldNeighbors, NewNeighbors) of
141     {trigger_rebuild, NewRT} ->
142         % trigger immediate rebuild
143         NewTriggerState = trigger:now(TriggerState),
144         ?RT:check(OldRT, NewRT, OldNeighbors, NewNeighbors, true),
145         new_state(NewNeighbors, NewRT, NewTriggerState);
146     {ok, NewRT} ->
147         ?RT:check(OldRT, NewRT, OldNeighbors, NewNeighbors, true),
148         new_state(NewNeighbors, NewRT, TriggerState)
149     end;
```

Every time a node's neighborhood changes, the `dht_node` sends an `update_rt` message to the routing table which will call `?RT:update/7` that decides whether the routing table should be rebuild. If so, it will stop any waiting trigger and schedule an immediate (periodic) stabilization.

### 9.3.2. Simple routing table (rt\_simple)

One implementation of a routing table is the `rt_simple`, which routes via the successor. Note that this is inefficient as it needs a linear number of hops to reach its goal. A more robust implementation, would use a successor list. This implementation is also not very efficient in the presence of churn.

## Data types

First, the data structure of the routing table is defined:

File `rt_simple.erl`:

```
26 -type key_t() :: 0..16#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF. % 128 bit numbers
27 -type rt_t() :: Succ::node:node_type().
28 -type external_rt_t() :: Succ::node:node_type().
29 -type custom_message() :: none().
```

The routing table only consists of a node (the successor). Keys in the overlay are identified by integers  $\geq 0$ .

## A simple `rm_beh` behaviour

File `rt_simple.erl`:

```
41 %% @doc Creates an "empty" routing table containing the successor.
42 empty(Neighbors) -> nodelist:succ(Neighbors).
```

File `rt_simple.erl`:

```
204 empty_ext(Neighbors) -> empty(Neighbors).
```

The empty routing table (internal or external) consists of the successor.

File `rt_simple.erl`:

Keys are hashed using MD5 and have a length of 128 bits.

File `rt_simple.erl`:

```
58 %% @doc Generates a random node id, i.e. a random 128-bit number.
59 get_random_node_id() ->
60     case config:read(key_creator) of
61     random -> hash_key_(randoms:getRandomId());
62     random_with_bit_mask ->
63         {Mask1, Mask2} = config:read(key_creator_bitmask),
64         (hash_key_(randoms:getRandomId()) band Mask2) bor Mask1
65     end.
```

Random node id generation uses the helpers provided by the `randoms` module.

File `rt_simple.erl`:

```
208 %% @doc Returns the next hop to contact for a lookup.
209 next_hop(State, _Key) -> node:pidX(dht_node_state:get(State, rt)).
```

Next hop is always the successor.

File `rt_simple.erl`:

```
73 %% @doc Triggered by a new stabilization round, renews the routing table.
74 init_stabilize(Neighbors, _RT) -> empty(Neighbors).
```

`init_stabilize/2` resets its routing table to the current successor.

File rt\_simple.erl:

```
78 %% @doc Updates the routing table due to a changed node ID, pred and/or succ.
79 -spec update(OldRT::rt(), OldNeighbors::odelist:neighborhood(),
80             NewNeighbors::odelist:neighborhood()) -> {ok, rt()}.
81 update(_OldRT, _OldNeighbors, NewNeighbors) ->
82     {ok, odelist:succ(NewNeighbors)}.
```

update/7 updates the routing table with the new successor.

File rt\_simple.erl:

```
86 %% @doc Removes dead nodes from the routing table (rely on periodic
87 %%      stabilization here).
88 filter_dead_node(RT, _DeadPid) -> RT.
```

filter\_dead\_node/2 does nothing, as only the successor is listed in the routing table and that is reset periodically in init\_stabilize/2.

File rt\_simple.erl:

```
92 %% @doc Returns the pids of the routing table entries.
93 to_pid_list(Succ) -> [node:pidX(Succ)].
```

to\_pid\_list/1 returns the pid of the successor.

File rt\_simple.erl:

```
97 %% @doc Returns the size of the routing table.
98 get_size(_RT) -> 1.
```

The size of the routing table is always 1.

File rt\_simple.erl:

```
136 %% @doc Returns the replicas of the given key.
137 get_replica_keys(Key) ->
138     [Key,
139      Key bxor 16#40000000000000000000000000000000,
140      Key bxor 16#80000000000000000000000000000000,
141      Key bxor 16#C0000000000000000000000000000000
142     ].
```

This get\_replica\_keys/1 implements symmetric replication.

File rt\_simple.erl:

```
102 %% @doc Returns the size of the address space.
103 n() -> n_().
104 %% @doc Helper for n/0 to make dialyzer happy with internal use of n/0.
105 -spec n_() -> 16#10000000000000000000000000000000.
106 n_() -> 16#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF + 1.
```

There are  $2^{128}$  available keys.

File rt\_simple.erl:

```
146 %% @doc Dumps the RT state for output in the web interface.
147 dump(Succ) -> [{"0", lists:flatten(io_lib:format("~p", [Succ]))}].
```

dump/1 lists the successor.

File `rt_simple.erl`:

```
219 %% @doc Converts the (external) representation of the routing table to a list
220 %%     in the order of the fingers, i.e. first=succ, second=shortest finger,
221 %%     third=next longer finger,...
222 to_list(State) -> [dht_node_state:get(State, rt)].
```

`to_list/1` lists the successor from the external routing table state.

File `rt_simple.erl`:

```
213 %% @doc Converts the internal RT to the external RT used by the dht_node. Both
214 %%     are the same here.
215 export_rt_to_dht_node(RT, _Neighbors) -> RT.
```

`export_rt_to_dht_node/2` states that the external routing table is the same as the internal table.

File `rt_simple.erl`:

```
165 %% @doc There are no custom messages here.
166 -spec handle_custom_message
167     (custom_message() | any(), rt_loop:state_active()) -> unknown_event.
168 handle_custom_message(_Message, _State) -> unknown_event.
```

Custom messages could be send from a routing table process on one node to the routing table process on another node and are independent from any other implementation.

File `rt_simple.hrl`:

```
172 %% @doc Notifies the dht_node and failure detector if the routing table changed.
173 %%     Provided for convenience (see check/5).
174 check(OldRT, NewRT, Neighbors, ReportToFD) ->
175     check(OldRT, NewRT, Neighbors, Neighbors, ReportToFD).
176
177 %% @doc Notifies the dht_node if the (external) routing table changed.
178 %%     Also updates the failure detector if ReportToFD is set.
179 %%     Note: the external routing table only changes the internal RT has
180 %%     changed.
181 check(OldRT, NewRT, _OldNeighbors, NewNeighbors, ReportToFD) ->
182     case OldRT == NewRT of
183     true -> ok;
184     _ ->
185         Pid = pid_groups:get_my(dht_node),
186         RT_ext = export_rt_to_dht_node(NewRT, NewNeighbors),
187         comm:send_local(Pid, {rt_update, RT_ext}),
188         % update failure detector:
189         case ReportToFD of
190         true ->
191             NewPids = to_pid_list(NewRT),
192             OldPids = to_pid_list(OldRT),
193             fd:update_subscriptions(OldPids, NewPids);
194         _ -> ok
195         end
196     end.
```

Checks whether the routing table changed and in this case sends the `dht_node` an updated (external) routing table state. Optionally the failure detector is updated. This may not be necessary, e.g. if `check` is called after a crashed node has been reported by the failure detector (the failure detector already unsubscribes the node in this case).

### 9.3.3. Chord routing table (`rt_chord`)

The file `rt_chord.erl` implements Chord's routing.

## Data types

File `rt_chord.erl`:

```
26 -type key_t() :: 0..16#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF. % 128 bit numbers
27 -type rt_t() :: gb_tree().
28 -type external_rt_t() :: gb_tree().
29 -type index() :: {pos_integer(), non_neg_integer()}.
30 -opaque custom_message() ::
31     {rt_get_node, Source_PID::comm:mypid(), Index::index()} |
32     {rt_get_node_response, Index::index(), Node::node:node_type()}.
```

The routing table is a `gb_tree`. Identifiers in the ring are integers. Note that in Erlang integer can be of arbitrary precision. For Chord, the identifiers are in  $[0, 2^{128})$ , i.e. 128-bit strings.

## The `rm_beh` behaviour for Chord (excerpt)

File `rt_chord.erl`:

```
46 %% @doc Creates an empty routing table.
47 empty(_Neighbors) -> gb_trees:empty().
```

File `rt_chord.erl`:

```
282 empty_ext(_Neighbors) -> gb_trees:empty().
```

`empty/1` returns an empty `gb_tree`, same for `empty_ext/1`.

`rt_chord:hash_key/1`, `rt_chord:get_random_node_id/0`, `rt_chord:get_replica_keys/1` and `rt_chord:n/0` are implemented like their counterparts in `rt_simple.erl`.

File `rt_chord.erl`:

```
286 %% @doc Returns the next hop to contact for a lookup.
287 %%     If the routing table has less entries than the rt_size_use_neighbors
288 %%     config parameter, the neighborhood is also searched in order to find a
289 %%     proper next hop.
290 %%     Note, that this code will be called from the dht_node process and
291 %%     it will thus have an external_rt!
292 next_hop(State, Id) ->
293     Neighbors = dht_node_state:get(State, neighbors),
294     case intervals:in(Id, nodelist:succ_range(Neighbors)) of
295     true -> node:pidX(nodelist:succ(Neighbors));
296     _ ->
297         % check routing table:
298         RT = dht_node_state:get(State, rt),
299         RTSize = get_size(RT),
300         NodeRT = case util:gb_trees_largest_smaller_than(Id, RT) of
301             {value, _Key, N} ->
302                 N;
303             nil when RTSize == 0 ->
304                 nodelist:succ(Neighbors);
305             nil -> % forward to largest finger
306                 {_Key, N} = gb_trees:largest(RT),
307                 N
308         end,
309         FinalNode =
310             case RTSize < config:read(rt_size_use_neighbors) of
311             false -> NodeRT;
312             _ ->
313                 % check neighborhood:
314                 nodelist:largest_smaller_than(Neighbors, Id, NodeRT)
315             end,
```

```

316         node:pidX(FinalNode)
317     end.

```

If the (external) routing table contains at least one item, the next hop is retrieved from the `gb_tree`. It will be the node with the largest id that is smaller than the id we are looking for. If the routing table is empty, the successor is chosen. However, if we haven't found the key in our routing table, the next hop will be our largest finger, i.e. entry.

File `rt_chord.erl`:

```

76 %% @doc Starts the stabilization routine.
77 init_stabilize(Neighbors, RT) ->
78     % calculate the longest finger
79     Id = nodelist:nodeid(Neighbors),
80     Key = calculateKey(Id, first_index()),
81     % trigger a lookup for Key
82     api_dht_raw:unreliable_lookup(Key, {send_to_group_member, routing_table,
83                                         {rt_get_node, comm:this(), first_index()}}),
84     RT.

```

The routing table stabilization is triggered for the first index and then runs asynchronously, as we do not want to block the `rt_loop` to perform other request while recalculating the routing table.

We have to find the node responsible for the calculated finger and therefore perform a lookup for the node with a `rt_get_node` message, including a reference to ourselves as the reply-to address and the index to be set.

The lookup performs an overlay routing by passing the message until the responsible node is found. There, the message is delivered to the `routing_table` process. The remote node sends the requested information back directly. It includes a reference to itself in a `rt_get_node_response` message. Both messages are handled by `rt_chord:handle_custom_message/2`:

File `rt_chord.erl`:

```

225 %% @doc Chord reacts on 'rt_get_node_response' messages in response to its
226 %%      'rt_get_node' messages.
227 -spec handle_custom_message
228     (custom_message(), rt_loop:state_active()) -> rt_loop:state_active();
229     (any(), rt_loop:state_active()) -> unknown_event.
230 handle_custom_message({rt_get_node, Source_PID, Index}, State) ->
231     MyNode = nodelist:node(rt_loop:get_neighb(State)),
232     comm:send(Source_PID, {rt_get_node_response, Index, MyNode}, ?SEND_OPTIONS),
233     State;
234 handle_custom_message({rt_get_node_response, Index, Node}, State) ->
235     OldRT = rt_loop:get_rt(State),
236     Neighbors = rt_loop:get_neighb(State),
237     NewRT = stabilize(Neighbors, OldRT, Index, Node),
238     check(OldRT, NewRT, rt_loop:get_neighb(State), true),
239     rt_loop:set_rt(State, NewRT);
240 handle_custom_message(_Message, _State) ->
241     unknown_event.

```

File `rt_chord.erl`:

```

150 %% @doc Updates one entry in the routing table and triggers the next update.
151 -spec stabilize(Neighbors::nodelist:neighborhood(), OldRT::rt(),
152                 Index::index(), Node::node:node_type()) -> NewRT::rt().
153 stabilize(Neighbors, RT, Index, Node) ->
154     MyId = nodelist:nodeid(Neighbors),
155     Succ = nodelist:succ(Neighbors),
156     case (node:id(Succ) /= node:id(Node)) % reached succ?
157     andalso (not intervals:in(           % there should be nothing shorter
158                 node:id(Node),           % than succ
159                 nodelist:succ_range(Neighbors))) of

```



```

160         true ->
161             NewRT = gb_trees:enter(Index, Node, RT),
162             NextKey = calculateKey(MyId, next_index(Index)),
163             CurrentKey = calculateKey(MyId, Index),
164             case CurrentKey /= NextKey of
165                 true ->
166                     Msg = {rt_get_node, comm:this(), next_index(Index)},
167                     api_dht_raw:unreliable_lookup(
168                         NextKey, {send_to_group_member, routing_table, Msg});
169                 _ -> ok
170             end,
171             NewRT;
172         _ -> RT
173     end.

```

stabilize/5 assigns the received routing table entry and triggers the routing table stabilization for the the next shorter entry using the same mechanisms as described above.

If the shortest finger is the successor, then filling the routing table is stopped, as no further new entries would occur. It is not necessary, that Index reaches 1 to make that happen. If less than  $2^{128}$  nodes participate in the system, it may happen earlier.

File rt\_chord.erl:

```

177 %% @doc Updates the routing table due to a changed node ID, pred and/or succ.
178 -spec update(OldRT::rt(), OldNeighbors::nodelist:neighborhood(),
179             NewNeighbors::nodelist:neighborhood()) -> {trigger_rebuild, rt()}.
180 update(_OldRT, _OldNeighbors, NewNeighbors) ->
181     % to be on the safe side ...
182     {trigger_rebuild, empty(NewNeighbors)}.

```

Tells the rt\_loop process to rebuild the routing table starting with an empty (internal) routing table state.

File rt\_chord.erl:

```

88 %% @doc Removes dead nodes from the routing table.
89 filter_dead_node(RT, DeadPid) ->
90     DeadIndices = [Index || {Index, Node} <- gb_trees:to_list(RT),
91                           node:same_process(Node, DeadPid)],
92     lists:foldl(fun(Index, Tree) -> gb_trees:delete(Index, Tree) end,
93                 RT, DeadIndices).

```

filter\_dead\_node removes dead entries from the gb\_tree.

File rt\_chord.erl:

```

321 export_rt_to_dht_node(RT, Neighbors) ->
322     Id = nodelist:nodeid(Neighbors),
323     Pred = nodelist:pred(Neighbors),
324     Succ = nodelist:succ(Neighbors),
325     Tree = gb_trees:enter(node:id(Succ), Succ,
326                           gb_trees:enter(node:id(Pred), Pred, gb_trees:empty())),
327     util:gb_trees_foldl(fun (_K, V, Acc) ->
328                           % only store the ring id and the according node structure
329                           case node:id(V) == Id of
330                               true -> Acc;
331                               false -> gb_trees:enter(node:id(V), V, Acc)
332                           end
333                           end, Tree, RT).

```

export\_rt\_to\_dht\_node converts the internal gb\_tree structure based on indices into the external representation optimised for look-ups, i.e. a gb\_tree with node ids and the nodes themselves.

File `rt_chord.hrl`:

```
245 %% @doc Notifies the dht_node and failure detector if the routing table changed.
246 %%     Provided for convenience (see check/5).
247 check(OldRT, NewRT, Neighbors, ReportToFD) ->
248     check(OldRT, NewRT, Neighbors, Neighbors, ReportToFD).
249
250 %% @doc Notifies the dht_node if the (external) routing table changed.
251 %%     Also updates the failure detector if ReportToFD is set.
252 %%     Note: the external routing table also changes if the Pred or Succ
253 %%     change.
254 check(OldRT, NewRT, OldNeighbors, NewNeighbors, ReportToFD) ->
255     case OldRT == NewRT andalso
256         nodelist:pred(OldNeighbors) == nodelist:pred(NewNeighbors) andalso
257         nodelist:succ(OldNeighbors) == nodelist:succ(NewNeighbors) of
258     true -> ok;
259     - ->
260         Pid = pid_groups:get_my(dht_node),
261         RT_ext = export_rt_to_dht_node(NewRT, NewNeighbors),
262         case Pid of
263             failed -> ok;
264             - -> comm:send_local(Pid, {rt_update, RT_ext})
265         end,
266         % update failure detector:
267         case ReportToFD of
268             true ->
269                 NewPids = to_pid_list(NewRT),
270                 OldPids = to_pid_list(OldRT),
271                 fd:update_subscriptions(OldPids, NewPids);
272             - -> ok
273         end
274     end.
```

Checks whether the routing table changed and in this case sends the `dht_node` an updated (external) routing table state. Optionally the failure detector is updated. This may not be necessary, e.g. if `check` is called after a crashed node has been reported by the failure detector (the failure detector already unsubscribes the node in this case).

## 9.4. Local Datastore

## 9.5. Cyclon

## 9.6. Vivaldi Coordinates

## 9.7. Estimated Global Information (Gossiping)

## 9.8. Load Balancing

## 9.9. Broadcast Trees

## 10. Transactions in Scalaris

### 10.1. The Paxos Module

### 10.2. Transactions using Paxos Commit

### 10.3. Applying the Tx-Modules to replicated DHTs

Introduces transaction processing on top of a Overlay

## 11. How a node joins the system

*Description is based on SVN revision r1370.*

After starting a new Scalaris-System as described in Section 3.2.1 on page 13, ten additional local nodes can be started by typing `admin:add_nodes(10)` in the Erlang-Shell that the management server opened <sup>1</sup>.

```
scalaris> ./bin/firstnode.sh
[...]  
(firstnode@csr-pc9)1> admin:add_nodes(10)
```

In the following we will trace what this function does in order to add additional nodes to the system. The function `admin:add_nodes(pos_integer())` is defined as follows.

File `admin.erl`:

```
39 % @doc add new Scalaris nodes on the local node  
40 -spec add_node_at_id(?RT:key()) -> pid_groups:groupname() | {error, term()}.  
41 add_node_at_id(Id) ->  
42     add_node([{{dht_node, id}, Id}, {skip_psv_lb}]).  
43  
44 -spec add_node([tuple()]) -> pid_groups:groupname() | {error, term()}.  
45 add_node(Options) ->  
46     DhtNodeId = randoms:getRandomId(),  
47     Desc = util:sup_supervisor_desc(  
48         DhtNodeId, config:read(dht_node_sup), start_link,  
49         [[{my_sup_dht_node_id, DhtNodeId} | Options]]),  
50     case supervisor:start_child(main_sup, Desc) of  
51         {ok, _Child, Group} -> Group;  
52         {error, already_present} -> add_node(Options); % try again, different Id  
53         {error, {already_started, _}} -> add_node(Options); % try again, different Id  
54         {error, _Error} = X -> X  
55     end.  
56  
57 -spec add_nodes(non_neg_integer()) -> {[pid_groups:groupname()], [{error, term()}]}.  
58 add_nodes(0) -> [];  
59 add_nodes(Count) ->  
60     Results = [add_node([]) || _X <- lists:seq(1, Count)],  
61     lists:partition(fun(E) -> not is_tuple(E) end, Results).
```

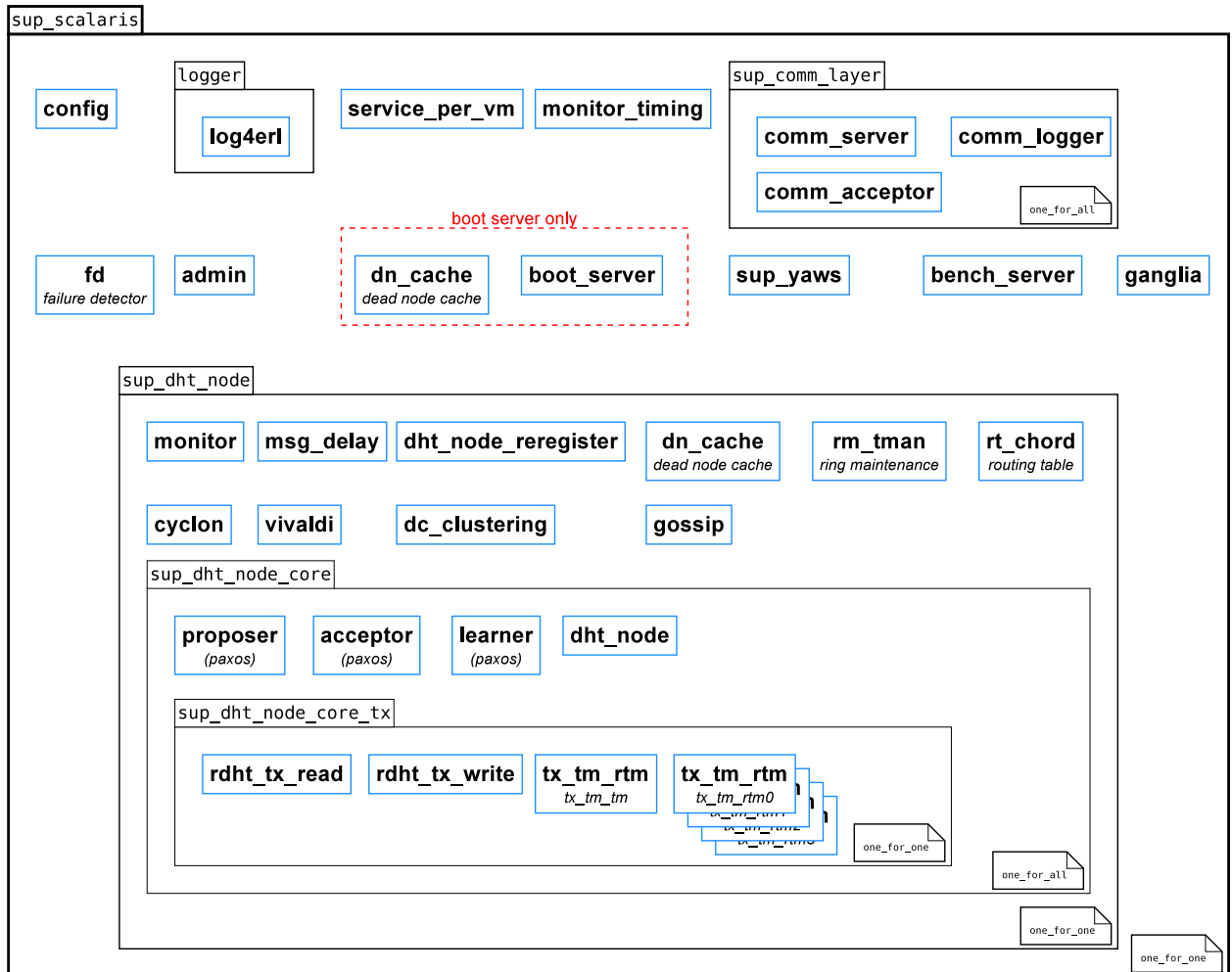
It calls `admin:add_node([])` Count times. This function starts a new child with the given options for the main supervisor `main_sup`. In particular, it sets a random ID that is passed to the new node as its suggested ID to join at. To actually perform the start, the function `sup_dht_node:start_link/1` is called by the Erlang supervisor mechanism. For more details on the OTP supervisor mechanism see Chapter 18 of the Erlang book [1] or the online documentation at <http://www.erlang.org/doc/man/supervisor.html>.

---

<sup>1</sup>Increase the log level to `info` to get more detailed startup logs. See Section 3.1.1 on page 12

## 11.1. Supervisor-tree of a Scalaris node

When a new Erlang VM with a Scalaris node is started, a `sup_scalaris` supervisor is started that creates further workers and supervisors according to the following scheme (processes starting order: left to right, top to bottom):



When new nodes are started using `admin:add_node/1`, only new `sup_dht_node` supervisors are started.

## 11.2. Starting the `sup_dht_node` supervisor and general processes of a node

Starting supervisors is a two step process: a call to `supervisor:start_link/2,3`, e.g. from a custom supervisor's own `start_link` method, will start the supervisor process. It will then call `Module:init/1` to find out about the restart strategy, maximum restart frequency and child processes. Note that `supervisor:start_link/2,3` will not return until `Module:init/1` has returned and all child processes have been started.

Let's have a look at `sup_dht_node:init/1`, the 'DHT node supervisor'.

File sup\_dht\_node.erl:

```
48 -spec init({pid_groups:groupname(), [tuple()]})
49       -> {ok, {{one_for_one, MaxRetries::pos_integer(), PeriodInSeconds::pos_integer()},
50               [ProcessDescr::any()]}}.
51 init({DHTNodeGroup, Options}) ->
52     pid_groups:join_as(DHTNodeGroup, ?MODULE),
53     mgmt_server:connect(),
54
55     Cyclon = util:sup_worker_desc(cyclon, cyclon, start_link, [DHTNodeGroup]),
56     DC_Clustering =
57         util:sup_worker_desc(dc_clustering, dc_clustering, start_link,
58                             [DHTNodeGroup]),
59     DeadNodeCache =
60         util:sup_worker_desc(deadnodecache, dn_cache, start_link,
61                             [DHTNodeGroup]),
62     Delayer =
63         util:sup_worker_desc(msg_delay, msg_delay, start_link,
64                             [DHTNodeGroup]),
65     Gossip =
66         util:sup_worker_desc(gossip, gossip, start_link, [DHTNodeGroup]),
67     Reregister =
68         util:sup_worker_desc(dht_node_reregister, dht_node_reregister,
69                             start_link, [DHTNodeGroup]),
70     RoutingTable =
71         util:sup_worker_desc(routing_table, rt_loop, start_link,
72                             [DHTNodeGroup]),
73     SupDHTNodeCore_AND =
74         util:sup_supervisor_desc(sup_dht_node_core, sup_dht_node_core,
75                                 start_link, [DHTNodeGroup, Options]),
76     Vivaldi =
77         util:sup_worker_desc(vivaldi, vivaldi, start_link, [DHTNodeGroup]),
78     Monitor =
79         util:sup_worker_desc(monitor, monitor, start_link, [DHTNodeGroup]),
80     MonitorPerf =
81         util:sup_worker_desc(monitor_perf, monitor_perf, start_link, [DHTNodeGroup]),
82     RepUpdate = case config:read(rep_update_activate) of
83                 true -> util:sup_worker_desc(rep_upd, rep_upd,
84                                             start_link, [DHTNodeGroup]);
85                 _ -> []
86             end,
87     %% order in the following list is the start order
88     {ok, {{one_for_one, 10, 1},
89         lists:flatten([
90             Monitor,
91             Delayer,
92             Reregister,
93             DeadNodeCache,
94             RoutingTable,
95             Cyclon,
96             Vivaldi,
97             DC_Clustering,
98             Gossip,
99             SupDHTNodeCore_AND,
100            MonitorPerf,
101            RepUpdate
102        ])}}.
```

The return value of the `init/1` function specifies the child processes of the supervisor and how to start them. Here, we define a list of processes to be observed by a `one_for_one` supervisor. The processes are: `Monitor`, `Delayer`, `Reregister`, `DeadNodeCache`, `RingMaintenance`, `RoutingTable`, `Cyclon`, `Vivaldi`, `DC_Clustering`, `Gossip` and a `SupDHTNodeCore_AND` process in this order.

The term `{one_for_one, 10, 1}` specifies that the supervisor should try 10 times to restart each process before giving up. `one_for_one` supervision means, that if a single process stops, only that process is restarted. The other processes run independently.

When the `sup_dht_node:init/1` is finished the supervisor module starts all the defined processes

by calling the functions that were defined in the returned list.

For a join of a new node, we are only interested in the starting of the SupDHTNodeCore\_AND process here. At that point in time, all other defined processes are already started and running.

### 11.3. Starting the sup\_dht\_node\_core supervisor with a peer and some paxos processes

Like any other supervisor the sup\_dht\_node\_core supervisor calls its sup\_dht\_node\_core:init/1 function:

File sup\_dht\_node\_core.erl:

```
40 -spec init({pid_groups:groupname(), Options::[tuple()]}) ->
41     {ok, {{one_for_all, MaxRetries::pos_integer(),
42             PeriodInSeconds::pos_integer(),
43             [ProcessDescr::any()]}},
44     init({DHTNodeGroup, Options}) ->
45     pid_groups:join_as(DHTNodeGroup, ?MODULE),
46     PaxosProcesses = util:sup_supervisor_desc(sup_paxos, sup_paxos,
47                                             start_link, [DHTNodeGroup, []]),
48     DHTNodeModule = config:read(dht_node),
49     DHTNode = util:sup_worker_desc(dht_node, DHTNodeModule, start_link,
50                                   [DHTNodeGroup, Options]),
51     DHTNodeMonitor = util:sup_worker_desc(
52         dht_node_monitor, dht_node_monitor, start_link,
53         [DHTNodeGroup, Options]),
54     TX =
55         util:sup_supervisor_desc(sup_dht_node_core_tx, sup_dht_node_core_tx, start_link,
56                                   [DHTNodeGroup]),
57     {ok, {{one_for_all, 10, 1},
58         [
59             PaxosProcesses,
60             DHTNodeMonitor,
61             DHTNode,
62             TX
63         ]}}.
```

It defines five processes, that have to be observed using a one\_for\_all-supervisor, which means, that if one fails, all have to be restarted. The dht\_node module implements the main component of a full Scalaris node which glues together all the other processes. Its dht\_node:start\_link/2 function will get the following parameters: (a) the processes' group that is used with the pid\_groups module and (b) a list of options for the dht\_node. The process group name was calculated a bit earlier in the code. *Exercise: Try to find where.*

File dht\_node.erl:

```
533 %% @doc spawns a scalaris node, called by the scalaris supervisor process
534 -spec start_link(pid_groups:groupname(), [tuple()]) -> {ok, pid()}.
535 start_link(DHTNodeGroup, Options) ->
536     gen_component:start_link(?MODULE, Options,
537                               [{pid_groups_join_as, DHTNodeGroup, dht_node}, wait_for_init]).
```

Like many other modules, the dht\_node module implements the gen\_component behaviour. This behaviour was developed by us to enable us to write code which is similar in syntax and semantics to the examples in [3]. Similar to the supervisor behaviour, a module implementing this behaviour has to provide an init/1 function, but here it is used to initialize the state of the component. This function is described in the next section.

Note: `?MODULE` is a predefined Erlang macro, which expands to the module name, the code belongs to (here: `dht_node`).

## 11.4. Initializing a `dht_node`-process

File `dht_node.erl`:

```
514 %% @doc joins this node in the ring and calls the main loop
515 -spec init(Options::[tuple()])
516     -> dht_node_state:state() |
517     {'$gen_component', [{on_handler, Handler::on_join}], State::dht_node_join:join_state()}.
518 init(Options) ->
519     {my_sup_dht_node_id, MySupDhtNode} = lists:keyfind(my_sup_dht_node_id, 1, Options),
520     erlang:put(my_sup_dht_node_id, MySupDhtNode),
521     % get my ID (if set, otherwise chose a random ID):
522     Id = case lists:keyfind({dht_node, id}, 1, Options) of
523         {{dht_node, id}, IdX} -> IdX;
524         _ -> ?RT:get_random_node_id()
525     end,
526     case is_first(Options) of
527         true -> dht_node_join:join_as_first(Id, 0, Options);
528         _ -> dht_node_join:join_as_other(Id, 0, Options)
529     end.
```

The `gen_component` behaviour registers the `dht_node` in the process dictionary. Formerly, the process had to do this itself, but we moved this code into the behaviour. If an ID was given to `dht_node:init/1` function as a `{{dht_node, id}, KEY}` tuple, the given `Id` will be used. Otherwise a random key is generated. Depending on whether the node is the first inside a VM marked as first or not, the according function in `dht_node_join` is called. Also the pid of the node's supervisor is kept for future reference.

## 11.5. Actually joining the ring

After retrieving its identifier, the node starts the join protocol which processes the appropriate messages calling `dht_node_join:process_join_state(Message, State)`. On the existing node, join messages will be processed by `dht_node_join:process_join_msg(Message, State)`.

### 11.5.1. A single node joining an empty ring

File `dht_node_join.erl`:

```
99 -spec join_as_first(Id::?RT:key(), IdVersion::non_neg_integer(), Options::[tuple()])
100     -> dht_node_state:state().
101 join_as_first(Id, IdVersion, _Options) ->
102     comm:init_and_wait_for_valid_pid(),
103     log:log(info, "[ Node ~w ] joining as first: (~.0p, ~.0p)",
104         [self(), Id, IdVersion]),
105     Me = node:new(comm:this(), Id, IdVersion),
106     % join complete, State is the first "State"
107     finish_join(Me, Me, Me, ?DB:new(), msg_queue:new()).
```

If the ring is empty, the joining node will be the only node in the ring and will thus be responsible for the whole key space. It will trigger all known nodes to initialize the comm layer and then finish the join. `dht_node_join:finish_join/5` just creates a new state for a Scalaris node consisting of



the given parameters (the node as itself, its predecessor and successor, an empty database and the queued messages that arrived during the join). It then activates all dependent processes and creates a routing table from this information.

The `dht_node_state:state()` type is defined in

File `dht_node_state.erl`:

```

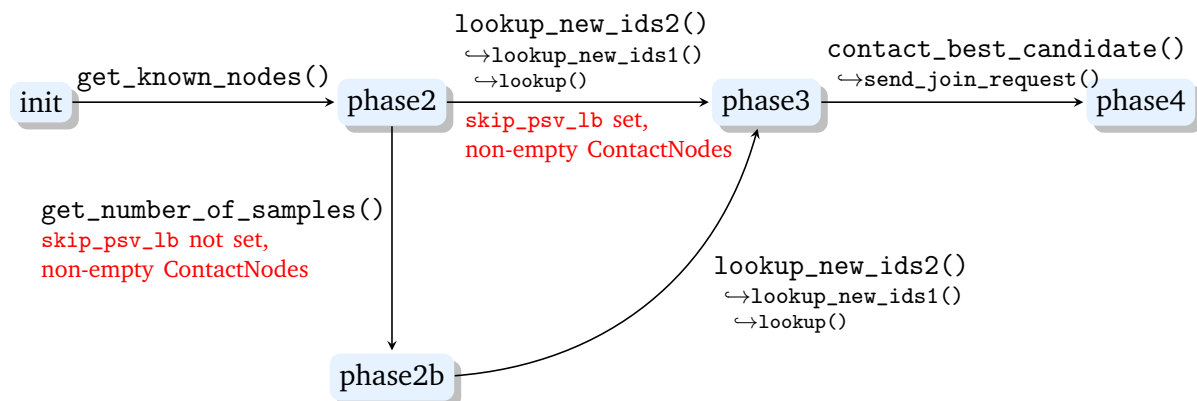
51 -record(state, {rt           = ?required(state, rt)           :: ?RT:external_rt(),
52                  rm_state    = ?required(state, rm_state)    :: rm_loop:state(),
53                  join_time    = ?required(state, join_time)   :: util:time(),
54                  db           = ?required(state, db)           :: ?DB:db(),
55                  tx_tp_db     = ?required(state, tx_tp_db)     :: any(),
56                  proposer     = ?required(state, proposer)    :: pid(),
57                  % slide with pred (must not overlap with 'slide with succ!'):
58                  slide_pred   = null :: slide_op:slide_op() | null,
59                  % slide with succ (must not overlap with 'slide with pred!'):
60                  slide_succ    = null :: slide_op:slide_op() | null,
61                  % additional range to respond to during a move:
62                  db_range     = [] :: [{intervals:interval(), slide_op:id()}],
63                  bulkowner_reply_timer = null :: null | reference(),
64                  bulkowner_reply_ids = [] :: [util:global_uid()],
65                  monitor_proc = ?required(state, monitor_proc) :: pid()
66                }).
67 -opaque state() :: #state{}.
```

### 11.5.2. A single node joining an existing (non-empty) ring

If a node joins an existing ring, its join protocol will step through the following four phases:

- **phase2** finding nodes to contact with the help of the configured `known_hosts`
- **phase2b** getting the number of Ids to sample (may be skipped)
- **phase3** lookup nodes responsible for all sampled Ids
- **phase4** joining a selected node and setting up item movements

The following figure shows a (non-exhaustive) overview of the transitions between the phases in the normal case. We will go through these step by step and discuss what happens if errors occur.



At first all nodes set in the `known_hosts` configuration parameter are contacted. Their responses are then handled in phase 2. In order to separate the join state from the ordinary `dht_node` state, the `gen_component` is instructed to use the `dht_node:on_join/2` message handler which delegates every message to `dht_node_join:process_join_state/2`.

File `dht_node_join.erl`:

```

111 -spec join_as_other(Id::?RT:key(), IdVersion::non_neg_integer(), Options::[tuple()])
112     -> {'$gen_component', [{on_handler, Handler::on_join}],
113         State::{join, phase2(), msg_queue:msg_queue()}}.
114 join_as_other(Id, IdVersion, Options) ->
115     comm:init_and_wait_for_valid_pid(),
116     log:log(info, "[ Node ~w ] joining, trying ID: (~.0p, ~.0p)",
117         [self(), Id, IdVersion]),
118     JoinUUID = util:get_pids_uid(),
119     get_known_nodes(JoinUUID),
120     msg_delay:send_local(get_join_timeout() div 1000, self(),
121         {join, timeout, JoinUUID}),
122     gen_component:change_handler(
123         {join, {phase2, JoinUUID, Options, IdVersion, [], [Id], []},
124             msg_queue:new()},
125         on_join).

```

## Phase 2 and 2b

Phase 2 collects all `dht_node` processes inside the contacted VMs. It therefore mainly processes `get_dht_nodes_response` messages and integrates all received nodes into the list of available connections. The next step depends on whether the `{skip_psv_lb}` option for skipping any passive load balancing algorithm has been given to the `dht_node` or not. If it is present, the node will only use the ID that has been initially passed to `dht_node_join:join_as_other/3`, issue a lookup for the responsible node and move to phase 3. Otherwise, the passive load balancing's `lb_psv_*:-get_number_of_samples/1` method will be called asking for the number of IDs to sample. Its answer will be processed in phase 2b.

`get_dht_nodes_response` messages arriving in phase 2b or later will be processed anyway and received `dht_node` processes will be integrated into the connections. These phases' operations will not be interrupted and nothing else is changed though.

File `dht_node_join.erl`:

```

153 % in phase 2 add the nodes and do lookups with them / get number of samples
154 process_join_state({get_dht_nodes_response, Nodes} = _Msg,
155     {join, JoinState, QueuedMessages})
156     when element(1, JoinState) == phase2 ->
157     ?TRACE_JOIN1(_Msg, JoinState),
158     Connections = [{null, Node} || Node <- Nodes, Node /= comm:this()],
159     JoinState1 = add_connections(Connections, JoinState, back),
160     NewJoinState = phase2_next_step(JoinState1, Connections),
161     ?TRACE_JOIN_STATE(NewJoinState),
162     {join, NewJoinState, QueuedMessages};
163
164 % in all other phases, just add the provided nodes:
165 process_join_state({get_dht_nodes_response, Nodes} = _Msg,
166     {join, JoinState, QueuedMessages})
167     when element(1, JoinState) == phase2b orelse
168         element(1, JoinState) == phase3 orelse
169         element(1, JoinState) == phase4 ->
170     ?TRACE_JOIN1(_Msg, JoinState),
171     Connections = [{null, Node} || Node <- Nodes, Node /= comm:this()],
172     JoinState1 = add_connections(Connections, JoinState, back),
173     ?TRACE_JOIN_STATE(JoinState1),
174     {join, JoinState1, QueuedMessages};

```

Phase 2b will handle `get_number_of_samples` messages from the passive load balance algorithm. Once received, new (unique) IDs will be sampled randomly so that the total number of join candidates (selected IDs together with fully processed candidates from further phases) is at least as high as the given number of samples. Afterwards, lookups will be created for all previous IDs as well as the new ones and the node will move to phase 3.

File dht\_node\_join.erl:

```
200 % note: although this message was send in phase2, also accept message in
201 % phase2, e.g. messages arriving from previous calls
202 process_join_state({join, get_number_of_samples, Samples, Conn} = _Msg,
203                   {join, JoinState, QueuedMessages})
204     when element(1, JoinState) =:= phase2 orelse
205         element(1, JoinState) =:= phase2b ->
206         ?TRACE_JOIN1(_Msg, JoinState),
207         % prefer node that send get_number_of_samples as first contact node
208         JoinState1 = reset_connection(Conn, JoinState),
209         % (re-)issue lookups for all existing IDs and
210         % create additional samples, if required
211         NewJoinState = lookup_new_ids2(Samples, JoinState1),
212         ?TRACE_JOIN_STATE(NewJoinState),
213         {join, NewJoinState, QueuedMessages};
214
215 % ignore message arriving in other phases:
216 process_join_state({join, get_number_of_samples, _Samples, Conn} = _Msg,
217                   {join, JoinState, QueuedMessages}) ->
218         ?TRACE_JOIN1(_Msg, JoinState),
219         NewJoinState = reset_connection(Conn, JoinState),
220         ?TRACE_JOIN_STATE(NewJoinState),
221         {join, NewJoinState, QueuedMessages};
```

Lookups will make Scalaris find the node currently responsible for a given ID and send a request to simulate a join to this node, i.e. a `get_candidate` message. Note that during such an operation, the joining node would become the existing node's predecessor. The simulation will be delegated to the passive load balance algorithm the joining node requested, as set by the `join_lb_psv` configuration parameter.

File dht\_node\_join.erl:

```
506 process_join_msg({join, get_candidate, Source_PID, Key, LbPsv, Conn} = _Msg, State) ->
507     ?TRACE1(_Msg, State),
508     LbPsv:create_join(State, Key, Source_PID, Conn);
```

## Phase 3

The result of the simulation will be send in a `get_candidate_response` message and will be processed in phase 3 of the joining node. It will be integrated into the list of processed candidates. If there are no more IDs left to process, the best among them will be contacted. Otherwise further `get_candidate_response` messages will be awaited. Such messages will also be processed in the other phases where the candidate will be simply added to the list.

File dht\_node\_join.erl:

```
253 process_join_state({join, get_candidate_response, OrigJoinId, Candidate, Conn} = _Msg,
254                   {join, JoinState, QueuedMessages})
255     when element(1, JoinState) =:= phase3 ->
256         ?TRACE_JOIN1(_Msg, JoinState),
257         JoinState0 = reset_connection(Conn, JoinState),
258         JoinState1 = remove_join_id(OrigJoinId, JoinState0),
259         JoinState2 = integrate_candidate(Candidate, JoinState1, front),
260         NewJoinState =
261             case get_join_ids(JoinState2) of
262             [] -> % no more join ids to look up -> join with the best:
263                 contact_best_candidate(JoinState2);
264             [_|_] -> % still some unprocessed join ids -> wait
265                 JoinState2
266             end,
267         ?TRACE_JOIN_STATE(NewJoinState),
268         {join, NewJoinState, QueuedMessages};
```

```

269
270 % In phase 2 or 2b, also add the candidate but do not continue.
271 % In phase 4, add the candidate to the end of the candidates as they are sorted
272 % and the join with the first has already started (use this candidate as backup
273 % if the join fails). Do not start a new join.
274 process_join_state({join, get_candidate_response, OrigJoinId, Candidate, Conn} = _Msg,
275                   {join, JoinState, QueuedMessages})
276   when element(1, JoinState) == phase2 orelse
277     element(1, JoinState) == phase2b orelse
278     element(1, JoinState) == phase4 ->
279     ?TRACE_JOIN1(_Msg, JoinState),
280     JoinState0 = reset_connection(Conn, JoinState),
281     JoinState1 = remove_join_id(OrigJoinId, JoinState0),
282     JoinState2 = case get_phase(JoinState1) of
283                   phase4 -> integrate_candidate(Candidate, JoinState1, back);
284                   _      -> integrate_candidate(Candidate, JoinState1, front)
285                 end,
286     ?TRACE_JOIN_STATE(JoinState2),
287     {join, JoinState2, QueuedMessages};

```

If `dht_node_join:contact_best_candidate/1` is called and candidates are available (there should be at this stage!), it will sort the candidates by using the passive load balance algorithm, send a `join_request` message and continue with phase 4.

File `dht_node_join.erl`:

```

801 %% @doc Contacts the best candidate among all stored candidates and sends a
802 %%      join_request (Timeouts = 0).
803 -spec contact_best_candidate(JoinState::phase_2_4())
804       -> phase2() | phase2b() | phase4().
805 contact_best_candidate(JoinState) ->
806   contact_best_candidate(JoinState, 0).
807 %% @doc Contacts the best candidate among all stored candidates and sends a
808 %%      join_request. Timeouts is the number of join_request_timeout messages
809 %%      previously received.
810 -spec contact_best_candidate(JoinState::phase_2_4(), Timeouts::non_neg_integer())
811       -> phase2() | phase2b() | phase4().
812 contact_best_candidate(JoinState, Timeouts) ->
813   JoinState1 = sort_candidates(JoinState),
814   send_join_request(JoinState1, Timeouts).

```

File `dht_node_join.erl`:

```

818 %% @doc Sends a join request to the first candidate. Timeouts is the number of
819 %%      join_request_timeout messages previously received.
820 %%      PreCond: the id has been set to the ID to join at and has been updated
821 %%      in JoinState.
822 -spec send_join_request(JoinState::phase_2_4(), Timeouts::non_neg_integer())
823       -> phase2() | phase2b() | phase4().
824 send_join_request(JoinState, Timeouts) ->
825   case get_candidates(JoinState) of
826   [] -> % no candidates -> start over (should not happen):
827     start_over(JoinState);
828   [BestCand | _] ->
829     Id = node_details:get(lb_op:get(BestCand, n1_new), new_key),
830     IdVersion = get_id_version(JoinState),
831     NewSucc = node_details:get(lb_op:get(BestCand, n1succ_new), node),
832     Me = node:new(comm:this(), Id, IdVersion),
833     CandId = lb_op:get(BestCand, id),
834     ?TRACE_SEND(node:pidX(NewSucc), {join, join_request, Me, CandId}),
835     comm:send(node:pidX(NewSucc), {join, join_request, Me, CandId}),
836     msg_delay:send_local(
837       get_join_request_timeout() div 1000, self(),
838       {join, join_request_timeout, Timeouts, CandId, get_join_uuid(JoinState)}),
839     set_phase(phase4, JoinState)
840   end.

```

The `join_request` message will be received by the existing node which will set up a slide operation with the new node. If it is not responsible for the key (anymore), it will deny the request and reply with a `{join, join_response, not_responsible, Node}` message. If it is responsible for the ID and is not participating in a slide with its current predecessor, it will set up a slide with the joining node:

File `dht_node_join.erl`:

```

525 process_join_msg({join, join_request, NewPred, CandId} = _Msg, State)
526   when (not is_atom(NewPred)) -> % avoid confusion with not_responsible message
527     ?TRACE1(_Msg, State),
528     TargetId = node:id(NewPred),
529     case dht_node_move:can_slide_pred(State, TargetId, {join, 'rcv'}) of
530       true ->
531         try
532           % TODO: implement step-wise join
533           MoveFullId = util:get_global_uid(),
534           Neighbors = dht_node_state:get(State, neighbors),
535           fd:subscribe([node:pidX(NewPred)], {move, MoveFullId}),
536           SlideOp = slide_op:new_sending_slide_join(
537             MoveFullId, NewPred, join, Neighbors),
538           SlideOp1 = slide_op:set_phase(SlideOp, wait_for_pred_update_join),
539           RMSubscrTag = {move, slide_op:get_id(SlideOp1)},
540           rm_loop:subscribe(self(), RMSubscrTag,
541             fun(_OldN, NewN, _IsSlide) ->
542               NewPred := nodelist:pred(NewN)
543             end,
544             fun dht_node_move:rm_notify_new_pred/4, 1),
545           State1 = dht_node_state:add_db_range(
546             State, slide_op:get_interval(SlideOp1),
547             slide_op:get_id(SlideOp1)),
548           MoveFullId = slide_op:get_id(SlideOp1),
549           MyOldPred = dht_node_state:get(State1, pred),
550           MyNode = dht_node_state:get(State1, node),
551           % no need to tell the ring maintenance -> the other node will trigger an update
552           % also this is better in case the other node dies during the join
553           %% rm_loop:notify_new_pred(comm:this(), NewPred),
554           Msg = {join, join_response, MyNode, MyOldPred, MoveFullId, CandId},
555           dht_node_move:send2(State1, SlideOp1, Msg)
556         catch throw:not_responsible ->
557           ?TRACE_SEND(node:pidX(NewPred),
558             {join, join_response, not_responsible, CandId}),
559           comm:send(node:pidX(NewPred),
560             {join, join_response, not_responsible, CandId}),
561           State
562         end;
563       - ->
564         ?TRACE(" [ ~.Op ] ~n ignoring join_request from ~.Op due to a running slide ~n",
565           [self(), NewPred]),
566         State
567     end;

```

## Phase 4

The joining node will receive the `join_response` message in phase 4 of the join protocol. If everything is ok, it will notify its ring maintenance process that it enters the ring, start all required processes and join the slide operation set up by the existing node in order to receive some of its data.

If the join candidate's node is not responsible for the candidate's ID anymore or the candidate's ID already exists, the next candidate is contacted until no further candidates are available and the join protocol starts over using `dht_node_join:start_over/1`.

Note that the `join_response` message will actually be processed in any phase. Therefore, if mes-

sages arrive late, the join can be processed immediately and the rest of the join protocol does not need to be executed again.

File dht\_node\_join.erl:

```

326 process_join_state({join, join_response, not_responsible, CandId} = _Msg,
327                   {join, JoinState, QueuedMessages} = State)
328   when element(1, JoinState) == phase4 ->
329     ?TRACE_JOIN1(_Msg, JoinState),
330     % the node we contacted is not responsible for the selected key anymore
331     % -> try the next candidate, if the message is related to the current candidate
332     case get_candidates(JoinState) of
333       [] -> % no candidates -> should not happen in phase4!
334         log:log(error, "[ Node ~w ] empty candidate list in join phase 4, "
335                      "starting over", [self()]),
336         NewJoinState = start_over(JoinState),
337         ?TRACE_JOIN_STATE(NewJoinState),
338         {join, NewJoinState, QueuedMessages};
339       [Candidate | _Rest] ->
340         case lb_op:get(Candidate, id) == CandId of
341           false -> State; % unrelated/old message
342           _ ->
343             log:log(info,
344                    "[ Node ~w ] node contacted for join is not responsible "
345                    "for the selected ID (anymore), trying next candidate",
346                    [self()]),
347             NewJoinState = try_next_candidate(JoinState),
348             ?TRACE_JOIN_STATE(NewJoinState),
349             {join, NewJoinState, QueuedMessages}
350         end
351     end;
352
353 % in other phases remove the candidate from the list (if it still exists):
354 process_join_state({join, join_response, not_responsible, CandId} = _Msg,
355                   {join, JoinState, QueuedMessages}) ->
356   ?TRACE_JOIN1(_Msg, JoinState),
357   {join, remove_candidate(CandId, JoinState), QueuedMessages};
358
359 % note: accept (delayed) join_response messages in any phase
360 process_join_state({join, join_response, Succ, Pred, MoveId, CandId} = _Msg,
361                   {join, JoinState, QueuedMessages} = State) ->
362   ?TRACE_JOIN1(_Msg, JoinState),
363   % only act on related messages, i.e. messages from the current candidate
364   Phase = get_phase(JoinState),
365   State1 = case get_candidates(JoinState) of
366     [] when Phase == phase4 -> % no candidates -> should not happen in phase4!
367       log:log(error, "[ Node ~w ] empty candidate list in join phase 4, "
368                    "starting over", [self()]),
369       NewJoinState = start_over(JoinState),
370       ?TRACE_JOIN_STATE(NewJoinState),
371       {join, NewJoinState, QueuedMessages};
372     [] -> State; % in all other phases, ignore the delayed join_response
373     % if no candidates exist
374     [Candidate | _Rest] ->
375       CandidateNode = node_details:get(lb_op:get(Candidate, n1succ_new), node),
376       CandidateNodeSame = node:same_process(CandidateNode, Succ),
377       case lb_op:get(Candidate, id) == CandId of
378         false ->
379           log:log(warn, "[ Node ~w ] ignoring old or unrelated "
380                        "join_response message", [self()]),
381           State; % ignore old/unrelated message
382         _ when not CandidateNodeSame ->
383           % id is correct but the node is not (should never happen!)
384           log:log(error, "[ Node ~w ] got join_response but the node "
385                        "changed, trying next candidate", [self()]),
386           NewJoinState = try_next_candidate(JoinState),
387           ?TRACE_JOIN_STATE(NewJoinState),
388           {join, NewJoinState, QueuedMessages};
389         _ ->
390           MyId = node_details:get(lb_op:get(Candidate, n1_new), new_key),
391           MyIdVersion = get_id_version(JoinState),

```

```

392         case MyId ==: node:id(Succ) orelse MyId ==: node:id(Pred) of
393         true ->
394             log:log(warn, "[ Node ~w ] chosen ID already exists, "
395                     "trying next candidate", [self()]),
396             % note: can not keep Id, even if skip_psv_lb is set
397             JoinState1 = remove_candidate_front(JoinState),
398             NewJoinState = contact_best_candidate(JoinState1),
399             ?TRACE_JOIN_STATE(NewJoinState),
400             {join, NewJoinState, QueuedMessages};
401         - ->
402             ?TRACE("[ ~.0p ]~n joined MyId:~.0p, MyIdVersion:~.0p~n "
403                   "Succ: ~.0p~n Pred: ~.0p~n",
404                   [self(), MyId, MyIdVersion, Succ, Pred]),
405             Me = node:new(comm:this(), MyId, MyIdVersion),
406             log:log(info, "[ Node ~w ] joined between ~w and ~w",
407                     [self(), Pred, Succ]),
408             rm_loop:notify_new_succ(node:pidX(Pred), Me),
409             rm_loop:notify_new_pred(node:pidX(Succ), Me),
410
411             finish_join_and_slide(Me, Pred, Succ, ?DB:new(),
412                                   QueuedMessages, MoveId)
413         end
414     end
415 end,
416 State1;

```

File dht\_node\_join.erl:

```

875 %% @doc Finishes the join and sends all queued messages.
876 -spec finish_join(Me::node:node_type(), Pred::node:node_type(),
877                 Succ::node:node_type(), DB::?DB:db(),
878                 QueuedMessages::msg_queue:msg_queue())
879         -> dht_node_state:state().
880 finish_join(Me, Pred, Succ, DB, QueuedMessages) ->
881     RMState = rm_loop:init(Me, Pred, Succ),
882     Neighbors = rm_loop:get_neighbors(RMState),
883     % wait for the ring maintenance to initialize and tell us its table ID
884     rt_loop:activate(Neighbors),
885     cyclon:activate(),
886     vivaldi:activate(),
887     dc_clustering:activate(),
888     gossip:activate(nodelist:node_range(Neighbors)),
889     dht_node_reregister:activate(),
890     msg_queue:send(QueuedMessages),
891     NewRT_ext = ?RT:empty_ext(Neighbors),
892     dht_node_state:new(NewRT_ext, RMState, DB).
893
894 %% @doc Finishes the join by setting up a slide operation to get the data from
895 %% the other node and sends all queued messages.
896 -spec finish_join_and_slide(Me::node:node_type(), Pred::node:node_type(),
897                            Succ::node:node_type(), DB::?DB:db(),
898                            QueuedMessages::msg_queue:msg_queue(), MoveId::slide_op:id())
899         -> {'$gen_component', [{on_handler, Handler::on}],
900             State::dht_node_state:state()}.
901 finish_join_and_slide(Me, Pred, Succ, DB, QueuedMessages, MoveId) ->
902     State = finish_join(Me, Pred, Succ, DB, QueuedMessages),
903     fd:subscribe([node:pidX(Succ)], {move, MoveId}),
904     SlideOp = slide_op:new_receiving_slide_join(MoveId, Pred, Succ, node:id(Me), join),
905     SlideOp1 = slide_op:set_phase(SlideOp, wait_for_node_update),
906     SlideOp2 = slide_op:set_msg_fwd(SlideOp1, slide_op:get_interval(SlideOp1)),
907     State1 = dht_node_state:set_slide(State, succ, SlideOp2),
908     RMSubscrTag = {move, slide_op:get_id(SlideOp2)},
909     comm:send_local(self(), {move, node_update, RMSubscrTag}),
910     gen_component:change_handler(State1, on).

```

The macro ?RT maps to the configured routing algorithm. It is defined in include/scalaris.hrl. For further details on the routing see Chapter 9.3 on page 40.



## Timeouts and other errors

The following table summarizes the timeout messages send during the join protocol on the joining node. It shows in which of the phases each of the messages is processed and describes (in short) what actions are taken. All of these messages are influenced by their respective config parameters, e.g. `join_timeout` parameter in the config files defines an overall timeout for the whole join operation. If it takes longer than `join_timeout` ms, a `{join, timeout}` will be send and processed as given in this table.

	<code>known_hosts_timeout</code>	<code>get_number_of_samples_timeout</code>	<code>lookup_timeout</code>	<code>join_request_timeout</code>	<code>timeout</code>
<b>phase2</b>	get known nodes from configured VMs	ignore	ignore	ignore	
<b>phase2b</b>	ignore	remove contact node, re-start join → phase 2 or 2b	ignore	ignore	
<b>phase3</b>	ignore	ignore	remove contact node, lookup remaining IDs → phase 2 or 3	ignore	
<b>phase3b</b>	ignore	ignore	ignore	ignore	re-start join → phase 2 or 2b
<b>phase4</b>	ignore	ignore	ignore	timeouts < 3? <sup>2</sup> → contact candidate otherwise: remove candidate no candidates left? → phase 2 or 2b otherwise: → contact next one → phase 3b or 4	

On the existing node, there is only one timeout message which is part of the join protocol: the `join_response_timeout`. It will be send when a slide operation is set up and if the timeout hits before the next message exchange, it will increase the slide operation's number of timeouts. The slide will be aborted if at least `join_response_timeouts` timeouts have been received. This parameter is set in the config file.

### Misc. (all phases)

Note that join-related messages arriving in other phases than those handling them will be ignored. Any other messages during a `dht_node`'s join will be queued and re-send when the join is complete.

<sup>2</sup>set by the `join_request_timeouts` config parameter



## 12. Directory Structure of the Source Code

The directory tree of Scalaris is structured as follows:

<code>bin</code>	contains shell scripts needed to work with Scalaris (e.g. start the management server, start a node, ...)
<code>contrib</code>	necessary third party packages (yaws and log4erl)
<code>doc</code>	generated Erlang documentation
<code>docroot</code>	root directory of the node's webserver
<code>ebin</code>	the compiled Erlang code (beam files)
<code>java-api</code>	a Java API to Scalaris
<code>log</code>	log files
<code>src</code>	contains the Scalaris source code
<code>test</code>	unit tests for Scalaris
<code>user-dev-guide</code>	contains the sources for this document

## 13. Java API

For the Java API documentation, we refer the reader to the documentation generated by javadoc or doxygen. The following commands create the documentation:

```
%> cd java-api  
%> ant doc  
%> doxygen
```

The documentation can then be found in `java-api/doc/index.html` (javadoc) and `java-api/doc-doxygen/html/index.html` (doxygen).

The API is divided into four classes:

- `de.zib.scalariz.Transaction` for (multiple) operations inside a transaction
- `de.zib.scalariz.TransactionSingleOp` for single transactional operations
- `de.zib.scalariz.ReplicatedDHT` for non-transactional (inconsistent) access to the replicated DHT items, e.g. deleting items
- `de.zib.scalariz.PubSub` for topic-based publish/subscribe operations

# Bibliography

- [1] Joe Armstrong. *Programming Erlang: Software for a Concurrent World*. Pragmatic Programmers, ISBN: 978-1-9343560-0-5, July 2007
- [2] Frank Dabek, Russ Cox, Frans Kaashoek, Robert Morris. *Vivaldi: A Decentralized Network Coordinate System*. ACM SIGCOMM 2004.
- [3] Rachid Guerraoui and Luis Rodrigues. *Introduction to Reliable Distributed Programming*. Springer-Verlag, 2006.
- [4] Ion Stoica, Robert Morris, David Karger, M. Frans Kaashoek and Hari Balakrishnan. *Chord: A Scalable Peer-to-peer Lookup Service for Internet Applications*. ACM SIGCOMM 2001, San Deigo, CA, August 2001, pp. 149-160. [http://pdos.csail.mit.edu/papers/chord:sigcomm01/chord\\_sigcomm.pdf](http://pdos.csail.mit.edu/papers/chord:sigcomm01/chord_sigcomm.pdf)
- [5] Márk Jelasity, Alberto Montresor, Ozalp Babaoglu. *T-Man: Gossip-based fast overlay topology construction*. Computer Networks (CN) 53(13):2321-2339, 2009.
- [6] F. Schintke, A. Reinefeld, S. Haridi, T. Schütt. *Enhanced Paxos Commit for Transactions on DHTs*. 10th IEEE/ACM Int. Conf. on Cluster, Cloud and Grid Computing, pp. 448-454, May 2010.
- [7] Spyros Voulgaris, Daniela Gavidia, Maarten van Steen. *CYCLON: Inexpensive Membership Management for Unstructured P2P Overlays*. J. Network Syst. Manage. 13(2): 2005.
- [8] Márk Jelasity, Alberto Montresor, Ozalp Babaoglu. *Gossip-based aggregation in large dynamic networks*. ACM Trans. Comput. Syst. 23(3), 219-252 (2005).

# Index

?RT  
  next\_hop, 41  
  update, 43

admin  
  add\_node, 52, 53

api\_tx, 39

comm, 3, 30, 30  
  get\_msg\_tag, 35  
  send\_to\_group\_member, 34

cs\_api, 31

cyclon, 39

dht\_node, 42, 43, 46, 50, 55, 58  
  init, 56  
  on\_join, 57

dht\_node\_join, 56  
  contact\_best\_candidate, 60, 60  
  finish\_join, 56, 63  
  finish\_join\_and\_slide, 63  
  join\_as\_other, 58  
  process\_join\_msg, 56  
  process\_join\_state, 56, 57  
  send\_join\_request, 60  
  start\_over, 61

dht\_node\_state  
  state, 57

erlang  
  exit, 32, 33  
  now, 29  
  send\_after, 29

ets  
  i, 29

fd, 38

gen\_component, 3, 29, 30, 30–38  
  bp\_barrier, 35, 36  
  bp\_cont, 35, 36  
  bp\_del, 35, 36  
  bp\_set, 35  
  bp\_set\_cond, 35  
  bp\_step, 36, 37  
  change\_handler, 32, 33, 33, 34  
  get\_state, 32  
  kill, 32, 34  
  post\_op, 34  
  runnable, 36  
  sleep, 34  
  start, 32  
  start\_link, 32, 33

intervals  
  in, 41

lb\_psv\_\*  
  get\_number\_of\_samples, 58

monitor, 38, 38, 39  
  check\_report, 39  
  client\_monitor\_set\_value, 39  
  monitor\_set\_value, 39  
  proc\_check\_timeslot, 39  
  proc\_set\_value, 39

msg\_delay, 29

paxos\_SUITE, 35  
  step\_until\_decide, 36

pdb, 38

pid\_groups, 3, 30, 30, 32–34, 55

randoms, 44

rm\_beh, 44, 47

routing\_table, 48

rrd, 38, 38  
  add, 38  
  add\_now, 38  
  create, 38

rt\_beh, 40  
  check, 42  
  check\_config, 42  
  dump, 42  
  empty, 42  
  empty\_ext, 42  
  export\_rt\_to\_dht\_node, 42  
  filter\_dead\_node, 42

- get\_random\_node\_id, [42](#)
- get\_replica\_keys, [42](#)
- get\_size, [42](#)
- handle\_custom\_message, [42](#)
- hash\_key, [42](#)
- init\_stabilize, [42](#)
- n, [42](#)
- next\_hop, [42](#)
- to\_list, [42](#)
- to\_pid\_list, [42](#)
- update, [42](#)
- rt\_chord, [46](#)
  - empty, [47](#)
  - empty\_ext, [47](#)
  - export\_rt\_to\_dht\_node, [49](#)
  - filter\_dead\_node, [49](#)
  - get\_random\_node\_id, [47](#)
  - get\_replica\_keys, [47](#)
  - handle\_custom\_message, [48](#), [48](#)
  - hash\_key, [47](#)
  - init\_stabilize, [48](#)
  - n, [47](#)
  - next\_hop, [47](#)
  - stabilize, [48](#)
  - update, [49](#)
- rt\_loop, [42](#), [43](#), [43](#), [49](#)
- rt\_simple, [43](#)
  - dump, [45](#)
  - empty, [44](#)
  - empty\_ext, [44](#)
  - export\_rt\_to\_dht\_node, [46](#)
  - filter\_dead\_node, [45](#)
  - get\_random\_node\_id, [44](#)
  - get\_replica\_keys, [45](#)
  - get\_size, [45](#)
  - handle\_custom\_message, [46](#)
  - hash\_key, [44](#)
  - init\_stabilize, [44](#)
  - n, [45](#)
  - next\_hop, [44](#)
  - to\_list, [45](#)
  - to\_pid\_list, [45](#)
  - update, [44](#)
- sup\_dht\_node
  - init, [53](#), [54](#)
  - start\_link, [52](#)
- sup\_dht\_node\_core, [55](#)
- sup\_scalaris, [53](#)
- supervisor
  - start\_link, [53](#)
- timer
  - sleep, [32](#)
  - tc, [29](#)
- util
  - tc, [29](#)
- vivaldi, [35](#)
- vivaldi\_latency, [35](#)
- your\_gen\_component
  - init, [32](#), [34](#)
  - on, [32-34](#)