

Scalaris

Users and Developers Guide

Version 0.1

Florian Schintke, Thorsten Schütt

May 15, 2009

Copyright 2007-2008 Konrad-Zuse-Zentrum für Informationstechnik Berlin

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

I	Users Guide	5
1	Introduction	7
2	Download and Installation	9
2.1	Requirements	9
2.2	Download	9
2.2.1	Development Branch	9
2.2.2	Releases	9
2.3	Configuration	9
2.4	Build	10
2.4.1	Linux	10
2.4.2	Windows	10
2.4.3	Java-API	10
2.5	Running Scalaris	11
2.5.1	Running on a local machine	11
2.5.2	Running distributed	11
2.6	Installation	12
2.7	Logging	12
3	Using the system	13
3.1	JSON API	13
3.1.1	Deleting a key	15
3.2	Java command line interface	16
3.3	Java API	16
4	Testing the system	17
4.1	Running the unit tests	17
II	Developers Guide	19
5	How a node joins the system	21
5.1	General Erlang server loop	21
5.2	Starting additional local nodes after boot	21
5.2.1	Supervisor-tree of a Scalaris node	22
5.2.2	Starting the or-supervisor and general processes of a node	22
5.2.3	Starting the and-supervisor with a peer and its local database	24
5.2.4	Initializing a cs_node -process	24
5.2.5	Actually joining the ring	25
5.2.6	Beginning to serve requests	27
6	Routing and routing tables in the Overlay	29

6.1	Simple routing table	30
6.1.1	Data types	31
6.1.2	A simple routingtable behaviour	31
6.2	Chord routing table	32
6.2.1	Data types	32
6.2.2	The routingtable behaviour for Chord	33
7	Directory Structure of the Source Code	35
8	Java API	37

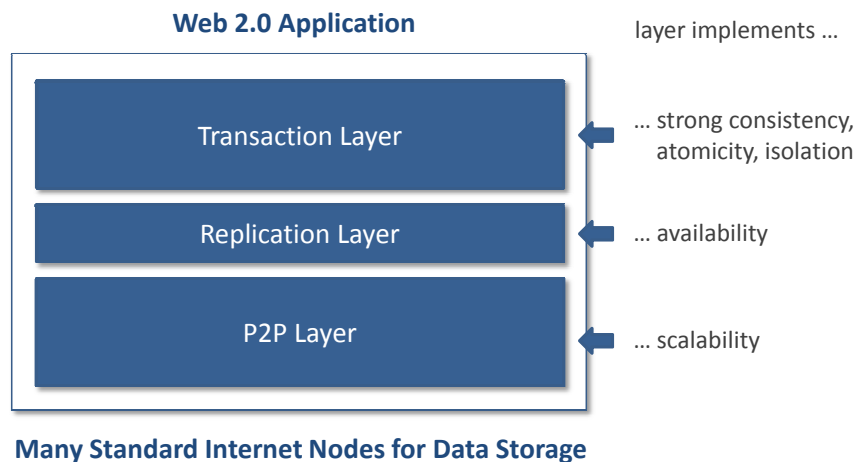
Part I

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 is partly funded by the EU projects Selfman and XtremOS. Additional information (papers, videos) can be found at <http://www.zib.de/CSR/Projects/scalaris> and <http://www.onscale.de/scalaris.html>.

2 Download and Installation

2.1 Requirements

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

- Erlang R12
- Erlang OTP (included in Erlang R12)
- GNU Make

Note, the Version 12 of Erlang is required. Scalaris will not work with older versions. To build the Java API the following modules are required additionally:

- Java Development Kit 1.6
- 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 1.6 installation, and **ANT_HOME** has to point to an Ant installation.

2.2 Download

The sources can be obtained from <http://code.google.com/p/scalaris>.

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 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 **bin/scalaris.local.cfg.example**. A local configuration file is necessary to run Scalaris on distributed nodes:

File **scalaris.local.cfg**:

```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
% Settings for distributed Erlang
```

```

% (see cs_send.erl to switch)

% {boot_host, {boot,'boot@foo.bar.com'}}.
% {log_host, {boot_logger, 'boot@foo.bar.com'}}.

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
% Settings for TCP mode.
% (see cs_send.erl to switch)

% Insert the appropriate IP-addresses for your setup
% as comma separated integers:
% IP Address, Port, and label of the boot server
{boot_host, {{127,0,0,1},14195,boot}}.

% IP Address, Port, and label of the log server
{log_host, {{127,0,0,1},14195,boot_logger}}.

```

boot_host defines the node where the boot server is running, which is contacted to join the system.

2.4 Build

2.4.1 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.

2.4.2 Windows

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

- Install Erlang
- Install OpenSSL (for crypto module)
- Checkout scalaris code from SVN
- copy an appropriate EMakefile_for from **contrib/win32** to the trunk-directory
- Adapt the path to your Erlang installation in build.bat
- run build.bat
- Go to the bin sub-directory
- Adapt the path to your Erlang installation in boot.bat
- run boot.bat

2.4.3 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
```

2.5 Running Scalaris

In Scalaris there are two kinds of processes:

- boot servers
- regular servers

In every Scalaris, at least one boot server is required. It will maintain a list of nodes taken part in the system and allows other nodes to join the ring. For redundancy, it is also possible to have several boot servers.

2.5.1 Running on a local machine

Open at least two shells. In the first, go into the `bin` directory:

```
%> cd bin
%> ./boot.sh
```

This will start the boot server. On success <http://localhost:8000> should point to the management interface page of the boot server. The main page will show you the number of nodes currently in the system. After a couple of seconds a first Scalaris should have started in the boot server and the number should increase to one. The main page will also allow you to store and retrieve key-value pairs.

In a second shell, you can now start a second Scalaris node. This will be a ‘regular server’. Go in the `bin` directory:

```
%> cd bin
%> ./cs_local.sh
```

The second node will read the configuration file and use this information to contact the boot server and will join the ring. 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:

```
%> cd bin
%> ./cs_local2.sh
```

In a fourth shell:

```
%> cd bin
%> ./cs_local3.sh
```

This will add 3 nodes to the network. The web pages at <http://localhost:8000> should show the additional nodes.

2.5.2 Running distributed

Scalaris can be installed on other machines in the same way as described in Sect. 2.6. In the default configuration, nodes will look for the boot server on localhost on port 14195. You should create a `scalaris.local.cfg` pointing to the node running the boot server.

```
% Insert the appropriate IP-addresses for your setup
% as comma separated integers:
% IP Address, Port, and label of the boot server
{boot_host, {{127,0,0,1},14195,boot}}.
```

If you are using the default configuration on the boot server it will listen on port 14195 and you only have to change the IP address in the configuration file. Otherwise the other nodes will not find the boot server. On the remote nodes, you only need to call `./cs_local.sh` and they will automatically contact the configured boot server.

2.6 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`. But is more convenient to build RPMs and install those.

```
svn checkout http://scalaris.googlecode.com/svn/trunk/ scalaris-0.0.1
tar -cvjf scalaris-0.0.1.tar.bz2 scalaris-0.0.1 --exclude-vcs
cp scalaris-0.0.1.tar.bz2 /usr/src/packages/SOURCES/
rpmbuild -ba scalaris-0.0.1/contrib/scalaris.spec
```

Your source and binary rpm will be generated in `/usr/src/packages/SRPMS` and `RPMS`. We also build rpms using checkouts from svn and provide them using the openSUSE BuildService at <http://download.opensuse.org/repositories/home:/tschuett/>. RPM packages are available for

- Fedora 9, 10,
- Mandriva 2008, 2009,
- openSUSE 11.0, 11.1,
- SLE 10, 11,
- CentOS 5 and
- RHEL 5.

Inside those repositories you will also find an erlang rpm - you don't need this if you already have a recent enough erlang version!

2.7 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`. The default value is `error`; only errors and severe problems are logged.

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

In some cases, it might be necessary to get more complete logging information, e.g. for debugging. In 5.2 on page 21, we are explaining the startup process of Scalarisnodes in more detail, here the `info` level provides more detailed information.

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

3 Using the system

3.1 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 over to a Scalaris node with a POST message. The result is an opaque TransLog and a list containing the results of the requests. 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 necessary. To finish the transaction, the request list can contain a 'commit' request as last element, which triggers the validation phase of the transaction processing.

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>). The following example illustrates the message flow:

Client

Make a transaction, that sets two keys:

```
{
  "method": "req_list",
  "version": "1.1",
  "params":
  [
    [
      { "write": { "keyA": "valueA" } },
      { "write": { "keyB": "valueB" } },
      { "commit": "commit" }
    ]
  ],
  "id": 0
}
```

Scalaris node

→

← Scalaris sends results back

```
{ "result":
  { "results":
    [
      { "op": "commit",
        "value": "ok",
        "key": "ok" },
      { "op": "write",
        "value": "valueB",
        "key": "keyB" },
      { "op": "write",
        "value": "valueA",
        "key": "keyA" }
    ],
    "translog":
    [...]
  },
  "id" : 0
}
```

In a second transaction: Read the two keys →

```
{
  "method": "req_list",
  "version": "1.1",
  "params":
    [
      [
        { "read": "keyA" },
        { "read": "keyB" }
      ]
    ]
  "id": 0
}
```

← Scalaris sends results back

```
{ "result":
  {"results":
    [
      { "op": "read",
        "value": "valueB",
        "key": "keyB" },
      { "op": "read",
        "value": "valueA",
        "key": "keyA" }
    ],
    "translog":
      [...] // this list is the translog
              // for further operations!
              // We name it TLOG here.
  },
  "id" : 0
}
```

Calculate something with the read values →
and make further requests, here a write and
the commit for the whole transaction. In-
clude also the latest translog we got from
Scalaris (named **TLOG** here).

```
{
  "method": "req_list",
  "version": "1.1",
  "params":
    [
      TLOG, // translog from prev. result.
      [
        { "write": { "keyA": "valueA2" } },
        { "commit": "commit" }
      ]
    ]
  "id" : 0
}
```

← Scalaris sends results back

```
{ "result":
  { "results":
    [ { "op": "commit",
        "value": "ok",
        "key": "ok" },
      { "op": "write",
        "value": "valueA2",
        "key": "keyA" }
    ],
    "translog":
    [...]
  },
  "id" : 0
}
```

A sample usage of the JSON API using Ruby can be found in `contrib/jsonrpc.rb`.

A single request list must not contain a key more than once!

The allowed requests are:

```
{ "read": "any_key" }

{ "write": { "any_key": "any_value" } }

{ "commit": "commit" }
```

The possible results are:

```
{ "op": "read", "key": "any_key", "value": "any_value" }
{ "op": "read", "key": "any_value", "fail": "reason" } // 'not_found' or 'timeout'

{ "op": "write", "key": "any_key", "value": "any_value" }
{ "op": "read", "key": "any_key", "fail": "reason" }

{ "op": "commit", "value": "ok", "key": "ok" }
{ "op": "commit", "value": "fail", "fail": "reason" }
```

3.1.1 Deleting a key

Outside transactions keys can also be deleted, but it has to be done with care, as explained in the following thread on the mailing list: http://groups.google.com/group/scalaris/browse_thread/thread/ff1d9237e218799.

```
{
  "method": "delete",
  "version": "1.1",
  "params":
  [
    { "key": "any_key" }
  ],
  "id" : 0
}
```

Two sample results

```
{ "result":
  { "ok": 2, // how many replicas were deleted successfully
    "results": [ "ok", "ok", "locks_set", "undef" ]
  }
}
```

```
{ "result":
  { "failure":"reason" }
}
```

3.2 Java command line interface

The jar file contains a small command line interface client. For convenience, we provide a wrapper script called **scalaris** which setups the Java environment:

```
%> cd java-api
%> ./scalaris -help
usage: scalaris
-g,--getsubscribers <topic>  get subscribers of a topic
-h                             print this message
-m,--minibench                run mini benchmark
-p,--publish <params>        publish a new message for a topic: <topic>
                               <message>
-r,--read <key>               read an item
-s,--subscribe <params>      subscribe to a topic: <topic> <url>
-u,--unsubscribe <params>    unsubscribe from a topic: <topic> <url>
-w,--write <params>          write an item: <key> <value>
```

Read and write can be used to read resp. write from/to the overlay. getsubscribers, publish, and subscribe are the PubSub functions.

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

The scalaris library requires that you are running a ‘regular server’ on the same node. Having a boot server running on the same node is not sufficient.

3.3 Java API

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

- **Scalaris** provides a high-level API similar to the command line client.
- **Transaction** provides a low-level API to the transaction mechanism.

For details we refer the reader to the Javadoc:

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


4 Testing the system

4.1 Running the unit tests

There are some unit tests in the **test** directory. 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 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 5 minutes. Only when the complete suite finished, it will present statistics on failed and successful tests.

Part II

Developers Guide

5 How a node joins the system

5.1 General Erlang server loop

Servers in Erlang often use the following structure to maintain a state while processing received messages:

```
receive
  Message ->
    State1 = f(State),
    loop(State1)
end.
```

The server runs an endless loop, that waits for a message, processes it and calls itself using tail-recursion in each branch. The loop works on a **State**, which can be modified when a message is handled.

5.2 Starting additional local nodes after boot

After booting a new Scalaris-System as described in Section 2.5.1 on page 11, ten additional local nodes can be started by typing `admin:add_nodes(10)` in the Erlang-Shell that the boot process opened ¹.

```
scalaris/bin> ./boot.sh
[...]
```

```
=INFO REPORT==== 12-May-2009::16:24:18 ===
Yaws: Listening to 0.0.0.0:8000 for servers
- http://localhost:8000 under ../docroot
[info] [ CC ] this() == {{127,0,0,1},14195}
[info] [ DNC <0.96.0> ] starting DeadNodeCache
[info] [ DNC <0.96.0> ] starting Dead Node Cache
[info] [ RM <0.97.0> ] starting ring maintainer

[info] [ RT <0.99.0> ] starting routintable
[info] [ Node <0.101.0> ] joining 315238232250031455306327244779560426902
[info] [ Node <0.101.0> ] join as first 315238232250031455306327244779560426902
[info] [ FD <0.74.0> ] starting pinger for {{127,0,0,1},14195,<0.101.0>}
[info] [ Node <0.101.0> ] joined
[info] [ CY ] Cyclon spawn: {{127,0,0,1},14195,<0.102.0>}
(boot@csr-pc9)1> admin:add_nodes(10)
```

In the following we will trace, what this function does to join additional nodes to the system. The function `admin:add_nodes(int)` is defined as follows.

File `admin.erl`:

```
37 %%-----
38 %% Function: add_nodes(int()) -> ok
39 %% Description: add new Scalaris nodes
40 %%-----
41 % @doc add new Scalaris nodes on the local node
42 % @spec add_nodes(int()) -> ok
43
44 add_nodes(Count) ->
```

¹Increase the log level to `info` to get the detailed startup logs. See Sect. 2.7 on page 12

```

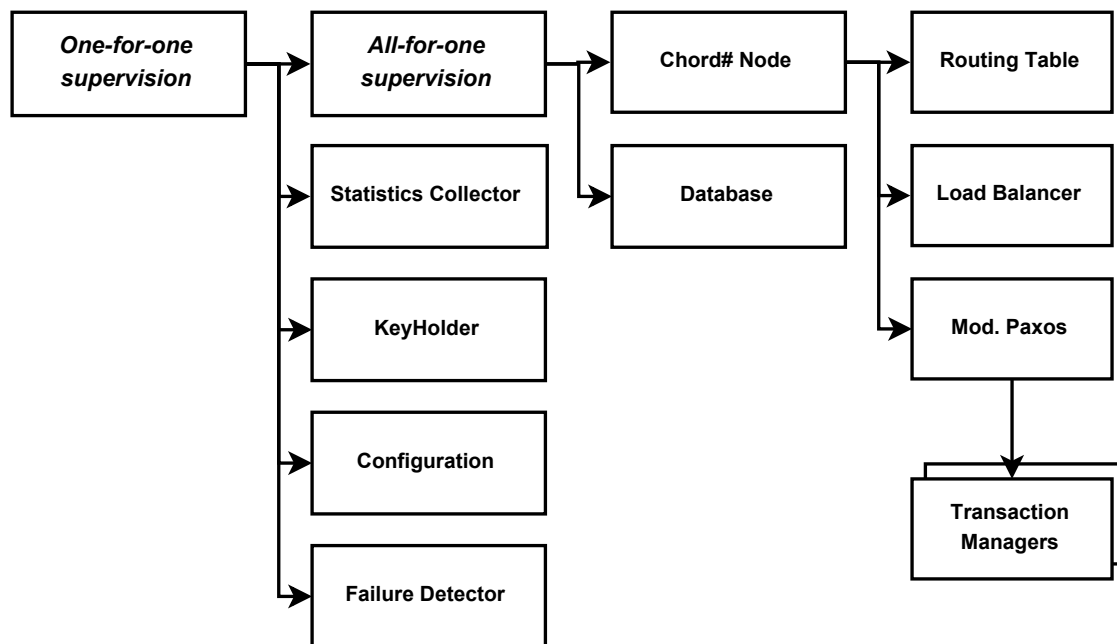
45     add_nodes(Count, 0).
46
47 % @spec add_nodes(int(), int()) -> ok
48 add_nodes(Count, Delay) ->
49     add_nodes_loop(Count, Delay).
50
51 add_nodes_loop(0, _) ->
52     ok;
53 add_nodes_loop(Count, Delay) ->
54     supervisor:start_child(main_sup, {randoms:getRandomId(),
55                                       {cs_sup_or, start_link, []},
56                                       permanent,
57                                       brutal_kill,
58                                       worker,
59                                       []}),
60     timer:sleep(Delay),
61     add_nodes_loop(Count - 1, Delay).

```

It calls `add_nodes_loop(Count, Delay)` with a delay of 0. This function starts a new child for the main supervisor `main_sup`. As defined by the parameters, to actually perform the start, the function `cs_sup_or:start_link` 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>.

5.2.1 Supervisor-tree of a Scalaris node

When starting a new node in the system, the following supervisor tree is created:



5.2.2 Starting the or-supervisor and general processes of a node

Starting supervisors is a two step process: the supervisor mechanism first calls the `init()` function of the defined module (`cs_sup_or:init()` in this case) and then calls the start function (`start_link` here).

So, let's have a look at `cs_sup_or:init`, the 'Scalaris *or* supervisor'.

File `cs_sup_or.erl`:

```

61 init([Options]) ->

```

```

62 InstanceId = string:concat("cs-node-", randoms:getRandomId()),
63 boot_server:connect(),
64 KeyHolder =
65     {cs_keyholder,
66      {cs_keyholder, start_link, [InstanceId]},
67      permanent,
68      brutal_kill,
69      worker,
70      []},
71 RSE =
72     {rse_chord,
73      {rse_chord, start_link, [InstanceId]},
74      permanent,
75      brutal_kill,
76      worker,
77      []},
78 Supervisor_AND =
79     {cs_supervisor_and,
80      {cs_sup_and, start_link, [InstanceId, Options]},
81      permanent,
82      brutal_kill,
83      supervisor,
84      []},
85 RingMaintenance =
86     {?RM,
87      {?RM, start_link, [InstanceId]},
88      permanent,
89      brutal_kill,
90      worker,
91      []},
92 RoutingTable =
93     {routingtable,
94      {rt_loop, start_link, [InstanceId]},
95      permanent,
96      brutal_kill,
97      worker,
98      []},
99 DeadNodeCache =
100     {deadnodecache,
101      {dn_cache, start_link, [InstanceId]},
102      permanent,
103      brutal_kill,
104      worker,
105      []},
106 {ok, {{one_for_one, 10, 1},
107      [
108          KeyHolder,
109          DeadNodeCache,
110          RingMaintenance,
111          RoutingTable,
112          Supervisor_AND
113          %RSE
114      ]}}}.

```

The return value of the `init()` 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: `KeyHolder`, `DeadNodeCache`, `RingMaintenance`, `RoutingTable`, and a `Supervisor_AND` process.

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.

The `cs_sup_or:init()` is finished and the supervisor module, starts all the defined processes by calling the functions that were defined in the list of the `cs_sup_or:init()`.

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

5.2.3 Starting the and-supervisor with a peer and its local database

Again, the OTP will first call the `init()` function of the corresponding module:

File `cs_sup_and.erl`:

```
58 init([InstanceId, Options]) ->
59     Node =
60         {cs_node,
61          {cs_node, start_link, [InstanceId, Options]},
62          permanent,
63          brutal_kill,
64          worker,
65          []},
66     DB =
67         {?DB,
68          {?DB, start_link, [InstanceId]},
69          permanent,
70          brutal_kill,
71          worker,
72          []},
73     Cyclon =
74         {cyclon,
75          {cyclon.cyclon, start_link, [InstanceId]},
76          permanent,
77          brutal_kill,
78          worker,
79          []},
80     {ok, {{one_for_all, 10, 1},
81          [
82              DB,
83              Node,
84              Cyclon
85          ]}}.
```

It defines three processes, that have to be observed using an `one_for_all`-supervisor, which means, that if one fails, all have to be restarted. Passed to the `init` function is the `InstanceId`, a random number to make nodes unique. It was calculated a bit earlier in the code. Exercise: Try to find where.

As you can see from the list, the `DB` is started before the `Node`. This is intended and important, because `cs_node` uses the database, but not vice versa. The supervisor first completely initializes the `DB` process and afterwards calls `cs_node:start_link`. We only go into details here, for the latter.

File `cs_node.erl`:

```
378 %% @doc spawns a scalaris node, called by the scalaris supervisor process
379 %% @spec start_link(term()) -> {ok, pid()}
380 start_link(InstanceId) ->
381     start_link(InstanceId, []).
382
383 start_link(InstanceId, Options) ->
384     gen_component:start_link(?MODULE, [InstanceId, Options], [{register, InstanceId, cs_node}]).
```

`cs_node` implements the `gen_component` behaviour. This component was developed by us to enable us to write code which is similar in syntax and semantics to the examples in [2]. Similar to the `supervisor` behaviour, the component has to provide an `init` 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: `cs_node`).

5.2.4 Initializing a `cs_node`-process

File `cs_node.erl`:

```
356 %% @doc joins this node in the ring and calls the main loop
357 -spec(init/1 :: ([any()]) -> cs_state:state()).
358 init([_InstanceId, Options]) ->
359     case lists:member(first, Options) of
360     true ->
361         ok;
362     false ->
363         timer:sleep(crypto:rand_uniform(1, 100) * 100)
364     end,
365     Id = cs_keyholder:get_key(),
366     {First, State} = cs_join:join(Id),
367     if
368         not First ->
369             cs_replica_stabilization:recreate_replicas(cs_state:get_my_range(State));
370     true ->
371         ok
372     end,
373     log:log(info, " [ Node ~w ] joined", [self()]),
374     State.
```

The `gen_component` behaviour registers the `cs_node` in the process dictionary. Formerly, the process had to do this himself, but we moved this code into the behaviour. If the `cs_node` is the first node, he will start immediately. Otherwise, the process sleeps for a random amount of time. If you would start 1000 processes with `admin:add_nodes(1000)`, the boot-server would receive many join requests at the same time, which is not intended. It will also make the ring stabilization process more complicated. Adding 100s of nodes within a short period of time induces more churn into the system, than the ring maintenance can handle.

Then, the node retrieves its `Id` from the keyholder: `Id = cs_keyholder:get_key()`. In the first call, a random identifier is returned, otherwise the latest set value. If the `cs_node`-process failed and is restarted by its supervisor, this call to the keyholder ensures, that the node still keeps its `Id`, assuming that the keyholder process is not failing. This is important for the load-balancing and for consistent responsibility of nodes to ensure consistent lookup in the structured overlay. Note: the name **Key-holder** actually is an id-holder.

If a node changes its position in the ring for load-balancing, the key-holder will be informed and the `cs_node` finishes itself. This triggers a restart of the corresponding database process via the and-supervisor. When the supervisor restarts both processes, they will retrieve the new position in the ring from the key-holder and join the ring there.

The supervisor was configured to restart a node at most 10 times. Does that mean, that a node can only change its position in the ring 10 times (caused by load-balancing)?

5.2.5 Actually joining the ring

After retrieving its identifier, the node starts the join process (`cs_join:join`).

File `cs_join.erl`:

```
87 %% @doc join a ring and return initial state
88 %%     the boolean indicates whether it was the first
89 %%     node in the ring or not
90 %% @spec join(Id) -> {true|false, state:state()}
91 %%     Id = term()
92 join(Id) ->
93     log:log(info, " [ Node ~w ] joining ~p", [self(), Id]),
94     Ringsize = boot_server:number_of_nodes(),
95     if
96         Ringsize == 0 ->
97             State = join_first(Id),
98             cs_reregister:reregister(),
```

```

99         {true, State};
100     true ->
101         case cs_lookup:reliable_get_node(erlang:get(instance_id),
102                                         Id, 60000) of
103             error ->
104                 join(Id);
105             {ok, Succ} ->
106                 State = join_ring(Id, Succ),
107                 cs_reregister:reregister(),
108                 {false, State}
109         end
110     end.

```

The boot-server is contacted to retrieve the known number of nodes in the ring. If the ring is empty, **join_first** is called. Otherwise, **join_ring** is called.

If the ring is empty, the joining node is the only node in the ring and will be responsible for the whole key space. **join_first** just creates a new state for a Scalaris node consisting of an empty routing table, a successorlist containing itself, itself as its predecessor, a reference to itself, its responsibility area from **Id** to **Id** (the full ring), and a load balancing schema.

File **cs_join.erl**:

```

50 %% @doc join an empty ring
51 join_first(Id) ->
52     log:log(info, "[ Node ~w ] join as first ~w", [self(), Id]),
53     Me = node:make(cs_send:this(), Id),
54     ?RM:initialize(Id, Me, Me, Me),
55     routingtable:initialize(Id, Me, Me),
56     cs_state:new(?RT:empty(Me), Me, Me, Me, {Id, Id}, cs_lb:new(), ?DB:new()).

```

The macro **?RT** maps to the configured routing algorithm and **?RM** to the configured ring maintenance algorithm. It is defined in **chordsharp.hrl**. For further details on the routing see Chapter 6 on page 29.

The state is defined in

File **cs_state.erl**:

```

57 new(RT, Successor, Predecessor, Me, MyRange, LB, DB) ->
58     #state{
59         routingtable = RT,
60         successor = Successor,
61         predecessor = Predecessor,
62         me = Me,
63         my_range = MyRange,
64         lb=LB,
65         join_time=now(),
66         deadnodes = gb_sets:new(),
67         trans_log = #translog{
68             tid_tm_mapping = dict:new(),
69             decided = gb_trees:empty(),
70             undecided = gb_trees:empty()
71         },
72         db = DB
73     }.

```

If a node joins an existing ring, **reliable_get_node** is called for the own **Id** in **cs_join:join()**. This lookup delivers the node who is currently responsible for the new node's identifier – the successor for the joining node. If this lookup fails for some reason, it is tried again, by recursively calling the **join()**.

What, if the **Id** is exactly the same as that of the existing node? This could lead to lookup and responsibility inconsistency? Can this be triggered by the load-balancing? This is a bug, that should be fixed!!!

Then, **cs_join:join_ring** is called:

File `cs_join.erl`:

```
61 join_ring(Id, Succ) ->
62     log:log(info, "[ Node ~w ] join_ring ~w", [self(), Id]),
63     Me = node:make(cs_send:this(), Id),
64     UniqueId = node:uniqueId(Me),
65     cs_send:send(node:pidX(Succ), {join, cs_send:this(), Id, UniqueId}),
66     receive
67         {join_response, Pred, Data} ->
68             log:log(info, "[ Node ~w ] got pred ~w", [self(), Pred]),
69             case node:is_null(Pred) of
70                 true ->
71                     DB = ?DB:add_data(?DB:new(), Data),
72                     ?RM:initialize(Id, Me, Pred, Succ),
73                     routingtable:initialize(Id, Pred, Succ),
74                     cs_state:new(?RT:empty(Succ), Succ, Pred, Me, {Id, Id}, cs_lb:new(), DB);
75                 false ->
76                     cs_send:send(node:pidX(Pred), {update_succ, Me}),
77                     DB = ?DB:add_data(?DB:new(), Data),
78                     ?RM:initialize(Id, Me, Pred, Succ),
79                     routingtable:initialize(Id, Pred, Succ),
80                     cs_state:new(?RT:empty(Succ), Succ, Pred, Me, {node:id(Pred), Id},
81                                 cs_lb:new(), DB)
82             end
83     end.
```

First the node is initialized. Then it sends a `join` message to the successor including a reference to itself and the chosen `Id`.

The message is received by the old node in `cs_node.erl`. There exists a `{join, X}` handler.

File `cs_node.erl`:

```
302 on({join, Source_PID, Id, UniqueId}, State) ->
303     cs_join:join_request(State, Source_PID, Id, UniqueId);
```

This triggers a call to `join_request` on the old node.

File `cs_join.erl`:

```
39 join_request(State, Source_PID, Id, UniqueId) ->
40     Pred = node:new(Source_PID, Id, UniqueId),
41     {DB, HisData} = ?DB:split_data(cs_state:get_db(State), cs_state:id(State), Id),
42     cs_send:send(Source_PID, {join_response, cs_state:pred(State), HisData}),
43     ?RM:update_pred(Pred),
44     cs_state:set_db(State, DB).
```

The `cs_node` notifies the ring maintenance, that he has a new predecessor. Then he removes the key-value pairs from his database which are now in the responsibility of the joining node. Then it sends a `join_response` to the new node with its former predecessor, the data, it has to host, and its successorlist.

Back on the joining node: it waits for the `join_response` message in `cs_join:join_ring()`. The next steps after the message was received from the old node are to initialize the maintenance components for the ring and routing table, the database and the state of the `cs_node`.

5.2.6 Beginning to serve requests

`cs_join:join()` was called from `cs_node:start()`, which now continues

File `cs_node.erl`:

```
356 %% @doc joins this node in the ring and calls the main loop
357 -spec(init/1 :: ([any()]) -> cs_state:state()).
358 init([_InstanceId, Options]) ->
359     case lists:member(first, Options) of
360         true ->
```

```

361         ok;
362         false ->
363             timer:sleep(crypto:rand_uniform(1, 100) * 100)
364     end,
365     Id = cs_keyholder:get_key(),
366     {First, State} = cs_join:join(Id),
367     if
368         not First ->
369             cs_replica_stabilization:recreate_replicas(cs_state:get_my_range(State));
370         true ->
371             ok
372     end,
373     log:log(info, "[ Node ~w ] joined", [self()],
374     State.

```

The `cs_replica_stabilization:recreate_replicas()` function is called, which is not yet implemented. It would recreated necessary replicas that were lost due to load-balancing and node failures.

Finally, the loop for request handling is started.

6 Routing and routing tables in the Overlay

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, so Scalaris provides a general interface to route a message to another peer, that is currently responsible for a given **key**.

File **cs_lookup.erl**:

```
[...]
unreliable_lookup(Key, Msg) ->
    get_pid(cs_node) ! {lookup_aux, Key, Msg}.

unreliable_get_key(Key) ->
    unreliable_lookup(Key, {get_key, cs_send:this(), Key}).
[...]
```

The message **Msg** could be a **get** 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 **cs_node.erl**.

The message routing is implemented in **lookup.erl**

File **lookup.erl**:

```
[...]
lookup_fin(Msg) ->
    self() ! Msg.

lookup_aux(State, Key, Msg) ->
    Terminate = util:is_between(cs_state:id(State), Key, cs_state:succ_id(State)),
    P = ?RT:next_hop(State, Key),
    ?LOG(" [ ~w | | Node    | ~w ] lookup_aux ~w ~w ~s~n",
        [calendar:universal_time(), self(), Terminate, P, Key]),
    if
        Terminate ->
            cs_send:send(P, {lookup_fin, Msg});
        true ->
            cs_send:send(P, {lookup_aux, Key, Msg})
    end.
[...]
```

Each node is responsible for a certain key interval. The function **util:is_between** is used to decide, whether the key is between the current node and its successor. If that is the case, final step is done using **lookup_fin()**, which delivers the 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**.

routingtable.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 ‘routingtable’.

The macro **?RT** is used to select the current implementation of routing tables. It is defined in **chordsharp.hrl**.

File **chordsharp.hrl**:

```
26 %%This file determines which kind of routingtable is used. Uncomment the
27 %%one that is desired.
```

```

28
29 %%Standard Chord routingtable
30 -define(RT, rt_chord).
31
32 %%Simple routingtable
33 -define(RT, rt_simple).

```

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

File `routingtable.erl`:

```

42 behaviour_info(callbacks) ->
43 [
44     % create a default routing table
45     {empty, 1},
46     % mapping: key space -> identifier space
47     {hash_key, 1}, {getRandomNodeId, 0},
48     % routing
49     {next_hop, 2},
50     % trigger for new stabilization round
51     {init_stabilize, 3},
52     % dead nodes filtering
53     {filterDeadNode, 2},
54     % statistics
55     {to_pid_list, 1}, {get_size, 1},
56     % for symmetric replication
57     {get_keys_for_replicas, 1},
58     % for debugging
59     {dump, 1},
60     % for bulkowner
61     {to_dict, 1}
62 ];

```

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

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

getRandomNodeId/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 routing table and a key and returns the node, that should be contacted next (is nearest to the id).

init_stabilize/3 is called periodically to rebuild the routing table. The parameters are the identifier of the node, the successor and the old routing table state.

filterDeadNode/2 is called by the faileddetector and tells the routing table about dead nodes to be eliminated from the routing table. This function cleans the routing table.

to_pid_list/1 get all PIDs of the routing table entries.

get_size/1 get the routing table's size.

get_keys_for_replicas/1 Returns for a given **Key** the keys of its replicas. This used for implementing symmetric replication.

dump/1 dump the state. Not mandatory, may just return **ok**.

to_dict/1 returns the routing tables entries in an array-like structure. This is used by bulk-operations to create a broadcast tree.

6.1 Simple routing table

One implementation of a routing table is the **rt_simple**, which routes via the successor, which 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 not very efficient on churn.

6.1.1 Data types

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

File `rt_simple.erl`:

```
39 % @type key(). Identifier.
40 -type(key()::pos_integer()).
41 % @type rt(). Routing Table.
42 -ifdef(types_are_builtin).
43 -type(rt()::{node:node_type(), gb_tree()}).
44 -else.
45 -type(rt()::{node:node_type(), gb_trees:gb_tree()}).
46 -endif.
```

A routing table is a pair of a node (the successor) and an (unused) `gb_tree`. Keys in the overlay are identified by integers.

6.1.2 A simple routingtable behaviour

File `rt_simple.erl`:

```
50 %% @doc creates an empty routing table.
51 %%     per default the empty routing should already include
52 %%     the successor
53 -spec(empty/1 :: (node:node_type() -> rt()).
54 empty(Succ) ->
55     {Succ, gb_trees:empty()}.
```

The empty routing table consists of the successor and an empty `gb_tree`.

File `rt_simple.erl`:

```
59 %% @doc hashes the key to the identifier space.
60 -spec(hash_key/1 :: (any() -> key()).
61 hash_key(Key) ->
62     BitString = binary_to_list(crypto:md5(Key)),
63     % binary to integer
64     lists:foldl(fun(E1, Total) -> (Total bsl 8) bor E1 end, 0, BitString).
```

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

File `rt_simple.erl`:

```
75 %% @doc returns the next hop to contact for a lookup
76 %% @spec next_hop(cs_state:state(), key()) -> pid()
77 next_hop(State, _Key) ->
78     cs_state:succ_pid(State).
```

Next hop is always the successor.

File `rt_simple.erl`:

```
82 %% @doc triggered by a new stabilization round
83 -spec(init_stabilize/3 :: (key(), node:node_type(), rt()) -> rt()).
84 init_stabilize(_Id, Succ, _RT) ->
85     % renew routing table
86     empty(Succ).
```

`init_stabilize/3` resets its routing table with the current successor.

File `rt_simple.erl`:

```
90 %% @doc removes dead nodes from the routing table
91 -spec(filterDeadNode/2 :: (rt(), cs_send:mypid()) -> rt()).
92 filterDeadNode(RT, _DeadPid) ->
93     RT.
```

`filterDeadNodes/2` does nothing, as only the successor is listed in the routing table and that is reset periodically in `init_stabilize/3`.

File `rt_simple.erl`:

```
97 %% @doc returns the pids of the routing table entries .
98 -spec(to_pid_list/1 :: (rt()) -> [cs_send:mypid()]).
99 to_pid_list({Succ, _RoutingTable} = _RT) ->
100     [node:pidX(Succ)].
```

`to_pid_list/1` returns the pids of the routing tables, as defined in `node.erl`.

File `rt_simple.erl`:

```
109 normalize(Key) ->
110     Key band 16#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF.
111
112 %% @doc returns the replicas of the given key
113 -spec(get_keys_for_replicas/1 :: (key() | string()) -> [key()]).
114 get_keys_for_replicas(Key) when is_integer(Key) ->
115     [Key,
116      normalize(Key + 16#40000000000000000000000000000000),
117      normalize(Key + 16#80000000000000000000000000000000),
118      normalize(Key + 16#C0000000000000000000000000000000)
119     ];
120 get_keys_for_replicas(Key) when is_list(Key) ->
121     get_keys_for_replicas(hash_key(Key)).
```

The `get_keys_for_replicas/1` implements symmetric replication, here. The call to `normalize` implements the modulo by throwing high bits away.

File `rt_simple.erl`:

```
126 %% @doc
127 -spec(dump/1 :: (rt()) -> ok).
128 dump(_State) ->
129     ok.
```

`dump/1` is not implemented.

6.2 Chord routing table

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

6.2.1 Data types

File `rt_chord.erl`:

```
40 -type(key() :: pos_integer()).
41 -ifdef(types_are_builtin).
42 -type(rt() :: gb_tree()).
43 -else.
44 -type(rt() :: gb_trees:gb_tree()).
45 -endif.
```

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.

6.2.2 The routingtable behaviour for Chord

File `rt_chord.erl`:

```
49 %% @doc creates an empty routing table.
50 -spec(empty/1 :: (node:node_type()) -> rt()).
51 empty(_Succ) ->
52     gb_trees:empty().
```

`empty/1` returns an empty `gb_tree`.

`hash_key(Key)` and `getRandomNodeId` call their counterparts from `rt_simple.erl`

File `rt_chord.erl`:

```
67 %% @doc returns the next hop to contact for a lookup
68 -spec(next_hop/2 :: (cs_state:state(), key()) -> cs_send:mypid()).
69 next_hop(State, Id) ->
70     case util:is_between(cs_state:id(State), Id, cs_state:succ_id(State)) of
71     %succ is responsible for the key
72     true ->
73         cs_state:succ_pid(State);
74     % check routing table
75     false ->
76         RT = cs_state:rt(State),
77         next_hop(cs_state:id(State), RT, Id, 127, cs_state:succ_pid(State))
78     end.
```

`next_hop` traverses the routing table beginning with the longest finger (2^{127}) by calling the helper function `next_hop/5`.

File `rt_chord.erl`:

```
82 % @private
83 -spec(next_hop/5 :: (key(), rt(), key(), pos_integer(), cs_send:mypid()) -> cs_send:mypid()).
84 next_hop(_N, _RT, _Id, 0, Candidate) -> Candidate;
85 next_hop(N, RT, Id, Index, Candidate) ->
86     case gb_trees:lookup(Index, RT) of
87     {value, Entry} ->
88         case util:is_between_closed(N, node:id(Entry), Id) of
89         true ->
90             node:pidX(Entry);
91         false ->
92             next_hop(N, RT, Id, Index - 1, Candidate)
93         end;
94     none ->
95         next_hop(N, RT, Id, Index - 1, Candidate)
96     end.
```

If the entry exists, it is retrieved from the `gb_tree`. If the id of the routing table entry is between ourselves and the searched id, the finger is chosen. If anything fails, **Candidate** (the successor) is chosen.

Why could a routing table entry be `null`? `filterDeadNodes` changes entries to `null`.

BUG: Instead of directly returning **Candidate** one should further traverse the routing table for shorter appropriate fingers. If doing so, a check whether **Index** is zero, would become necessary.

If the finger is too long, recursively try the next shorter finger.

File `rt_chord.erl`:

```
100 %% @doc starts the stabilization routine
101 -spec(init_stabilize/3 :: (key(), node:node_type(), rt()) -> rt()).
102 init_stabilize(Id, Succ, RT) ->
103     % calculate the longest finger
104     Key = calculateKey(Id, 127),
```

```

105 % trigger a lookup for Key
106 cs_lookup:unreliable_lookup(Key, {rt_get_node, cs_send:this(), 127}),
107 cleanup(gb_trees:iterator(RT), RT, Succ).

```

The routing table stabilization is triggered with the index 127 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 **cs_node**. At the destination the message is handled in **cs_node.erl**:

File **cs_node.erl**:

```

193 on({rt_get_node, Source_PID, Cookie}, State) ->
194   cs_send:send(Source_PID, {rt_get_node_response, Cookie, cs_state:me(State)}),
195   State;

```

The remote node just sends the requested information back directly in a **rt_get_node_response** message including a reference to itself. When receiving the routing table entry, we call **stabilize/5**.

File **rt_chord.erl**:

```

142 %% @doc updates one entry in the routing table
143 %% and triggers the next update
144 -spec(stabilize/5 :: (key(), node:node_type(), rt(), pos_integer(), node:node_type()) -> rt()).
145 stabilize(Id, Succ, RT, Index, Node) ->
146   case node:is_null(Node) of
147     true ->
148       RT;
149     false ->
150       case (node:id(Succ) == node:id(Node)) or (Id == node:id(Node)) or (Index == -1) of
151         true ->
152           % delete lower entries
153           prune_table(RT, Index);
154         false ->
155           NewRT = gb_trees:enter(Index, Node, RT),
156           Key = calculateKey(Id, Index - 1),
157           cs_lookup:unreliable_lookup(Key, {rt_get_node, cs_send:this(), Index - 1}),
158           NewRT
159       end
160   end.

```

stabilize/5 assigns the received routing table entry and triggers to fill the next shorter one using the same mechanisms as described.

When 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.

filterDeadNode removes dead entries from the **gb_tree**.

File **rt_chord.erl**:

```

111 %% @doc remove all entries
112 -spec(filterDeadNode/2 :: (rt(), cs_send:mypid()) -> rt()).
113 filterDeadNode(RT, DeadPid) ->
114   DeadIndices = [Index || {Index, Node} <- gb_trees:to_list(RT),
115                        node:pidX(Node) == DeadPid],
116   lists:foldl(fun (Index, Tree) -> gb_trees:delete(Index, Tree) end,
117               RT, DeadIndices).

```

7 Directory Structure of the Source Code

The directory tree of Scalaris is structured as follows:

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

8 Java API

For the Java API documentation, we refer the reader to Javadoc resp. doxygen. The following commands create the documentation:

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

The Javadoc can be found in **java-api/doc/index.html**. The doxygen is in **doc-doxygen/html/index.html**.

We provide two kinds of APIs:

- high-level access with **de.zib.scalarisc.Scalaris**
- low-level access with **de.zib.scalarisc.Transaction**

The former provides general functions for reading and writing single key-value pairs and an API for the built-in PubSub-service. The latter allows the user to write custom transactions which can modify an arbitrary number of key-value pairs within one transaction.

Bibliography

- [1] Joe Armstrong. *Programming Erlang: Software for a Concurrent World*. Pragmatic Programmers, ISBN: 978-1-9343560-0-5, July 2007
- [2] Rachid Guerraoui and Luis Rodrigues. *Introduction to Reliable Distributed Programming*. Springer-Verlag, 2006.