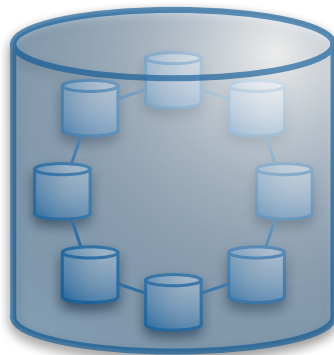# A Distributed Low-Level Data Structure for Composing Scalable Concurrent Web Services

**Master Thesis in Computer Science**

Submitted by
Michael Kussmaul
kussmaul@nix.ch
Zurich, Switzerland
Student ID: 96-716-055


**Department of Information Technology**
**University of Zurich, Switzerland**
**Prof. Dr. Klaus R. Dittrich**


Supervisor
Dr. Suo Cong

November 17, 2004

## Abstract

**English** - Distributed storage/retrieval of information is becoming increasingly crucial and is used extensively in various kinds of applications, such as web search engines, real-time systems, high performance computing, grid computing, and distributed file systems, to achieve higher throughput (enabling parallel access to the data), and more robust performance (redundancy in case of failures). While various alternatives have been proposed, this work will primarily focus on low-level mechanisms, and will introduce a new approach for a distributed key-value storage based on a b-tree implementation.

**Deutsch** - Verteiltes Speichern/Abfragen von Information wird je länger je mehr wichtiger und wird in verschiedensten Anwendungen benutzt. Beispiele sind Suchmaschinen, Echtzeitsysteme, Hochleistungsrechner, Grid-Computing und verteilte Dateisysteme, um höheren Durchsatz (ermöglicht parallelen Zugriff auf Daten) und Robustheit (Redundanz bei Ausfällen) zu ermöglichen. Während es verschiedene Möglichkeiten gibt, wird sich diese Arbeit auf grundlegende Mechanismen konzentrieren und einen neuen Ansatz einführen, der auf einem verteilten Assoziativspeicher (key-value) Mechanismus einer B-Baum Implementation basiert.
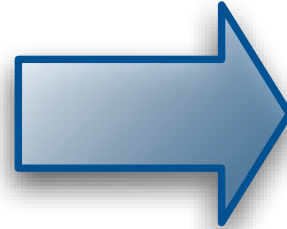
# Contents

# 1 Introduction

## 1.1 Situation

The increasing amount of data [18, 70] and the increasing level of access to those data has changed the way in which data is stored.

At the time when most basic systems (such as file-systems, databases and email systems) were designed, the amount of data they needed to store was not that great and everything could be handled by one single computer (e.g. a mainframe [81]). Over time, the amount of data has grown exponentially, and today's systems need to handle large amounts of data. Additionally, the data is stored in multiple places, and systems access this data from different sources. The number of users of such systems has also grown exponentially.

In today's world, a monolithic system, which only runs on a single machine, will not be able to cope with the amount of data and users accessing the system. Over time, most systems have moved from a centralized approach to a distributed one, which should allow for load-sharing the data and accessing requests from a broader range of machines.

It is interesting to note the similarities between those systems: The underlying storage mechanism is based on a key-value architecture for most of those systems. This means the system stores the data in chunks and stores each of them under a certain key. This key allows the user to retrieve the data again afterwards. The approach is similar to a conventional encyclopedia, where the keys are the items and the values are the definitions, as shown in Figure 1.

| Key | Value |
|---|---|
| *Zündapp* | *Was a famous German motorcycle brand.* |
| *Zürich* | *[ˈtsyːrɪç] Largest city in Switzerland (population 364,558 in 2002) and capital of the canton of Zürich.* |
| *Zug* | *City in central Switzerland (population 23,000 in 2004) and capital of the canton of Zug.* |

Figure 1: Key-Value Pairs in an Encyclopedia

## 1.2    Problems

There are several ways to actually distribute a system and most of the methods are highly dependent on the actual environment and usage case, e.g.:

- **File Systems** - Transformed from local file systems to network file systems (e.g. NFS [83]) and towards storage area networks, where the actual storage can be distributed to several machines.

- **Databases** - Transformed to a more distributed approach, e.g. by splitting the database tables into several chunks (partitioning). Depending on the architecture, they allow distribution to a few or several hundred machines.

While most of those systems already try to solve the problem of distributed storage, they have done the distribution on a high-level; the fundamental underlying structure has not been dramatically changed. This could soon be a limit, as the data continues to grow exponentially and a more radical and general approach for distributing the data could be necessary.

So much more radical an approach would be a low-level distribution already at the key-value layer. The system on top of it would then already be abstracted from the complexity of a distributed system, as the access to its data would still be the same key-value approach as before, with the difference that those key-values are now already distributed among multiple nodes. Also, more general frameworks could be implemented this way, by offloading all storage needs to a worldwide distributed storage network, so all large scale systems could access data from it, just by using a well-specified interface around a key-value principle.

Figure 2: Distributed key-value pairs form a worldwide storage network.

Currently, great efforts to solve this problem are being made in the research community. For the more relaxed case of unsorted key-value (hash tables) architectures, some algorithms and implementations already exist (they are called DHT: distributed hash tables) and are already being used in some special applications, such as file-sharing, peer-to-peer networks, but are not suited for other applications (like databases) where the sorted access to data is important.

For the more interesting case of distributed *sorted* key-value architectures, not much has been done so far. There are some research project working towards solving this problem (e.g. P-Rings [14], Boxwood [49]) but aside from some private prototypes, nothing "real" has been released to the public yet.

## 1.3    ERT (Extended Root Tree) Approach

The goal of this work is to build a prototype of a system that provides a distributed storage for *sorted* key-value pairs. It can then be used for different kind of applications (e.g. distributed databases, file systems, web applications).

ERT (Extended Root Tree) is a proposed solution for such a distributed sorted key-value architecture and will be introduced in detail in section 3. It is based on a conventional balanced tree (b-Tree) architecture, which has been extended by a separate handling of their root nodes (Extended Root Tree).

Figure 3: Extended Root Tree

This data structure, which provides distributed storage and a simple interface around its key-value architecture for small/medium/large systems, can be used as an underlying framework for building different kinds of applications on top of it.

The prototype implementation of ERT is hosted at `http://www.nix.ch/kussmaul/ert` and has been released under a BSD open source license [36].

## 1.4   Structure

The work can roughly be grouped into 3 main parts: The first part describes the theoretical background, the second part describes the chosen approach and the third part the implementation details. Additionally there is a section about "Evaluation" and an ending section, the "Conclusion".

- **Distributed Data Storage** - Describes state-of-the art methods for managing distributed data. It will first give an overview of general distribution techniques and describe the key factors for a successful distributed system. It will then describe methods for distributing data tailored to requirements for small/medium/large systems. The next part will describe local methods for storing data, which are already being used in monolithic systems and are still needed for distributed systems as well. The last part will then discuss performance aspects of a distributed system, such as concurrent access to data and high throughput systems.

- **Overview of ERT (Extended Root Tree)** - Describes the chosen approach to build such a distributed data structure on a more theoretical level.

- **Implementation of ERT (Extended Root Tree)** - Describes the current JAVA implementation of the built prototype and is meant as a program documentation to understand the underlying implementation details.

- **Evaluation** - Describes the benchmarking of the prototype and strengths and weaknesses found during testing and implementation.

- **Conclusion** - Provides a summary and a short outlook of the chosen approach.

# 2 Distributed Data Storage

## 2.1 Introduction

### 2.1.1 Nature of Data

Storage and retrieval of data is used everywhere and is becoming more and more important. In today's world we see the following tendencies:

- **Data Growth** - As we live in an information age, the amount of data (conventional, digital) is increasing rapidly. Studies [18, 70] have tried to get an overview of how much data is produced each year and estimate that in the year 2002 alone, more than 5 terabytes of data were produced [70]. This includes data stored on paper, film, magnetic and optical mediums. The biggest part of the data growth has come from magnetic storage. Digitization is also becoming more common: Old movies are being re-released in digital form (DVD, cinema), and books are being digitized (E.g., Google.com provides a service for full-text searching of books [29]).

- **Data Distribution** - Data needed for a specific task is often not centralized anymore; it is distributed and stored at several locations, even around the world. In the future, this tendency will continue to grow. As for certain data, specific requirements for storage are needed: High-value data (such as permanent data for banking institutes and telecommunications firms) is stored at special secured sites with access control and protection against power failure or other disasters. In contrast, data of lower value (such as temporary data for product downloads, web space) is more commonly stored at places, which provide the most value for the money. Some examples of distribution are:

  - **Personal Data** - More and more people already work on several computers at different locations, such as on a computer at work, at home and perhaps on a mobile device (such as laptop, palmtop or mobile phone). Some data reside on a specific computer, but there is a trend to "synchronize" the data among all devices, or to allow access to important data from all devices (such as emails, contacts, appointments).

- **Business Data** - Since companies have been growing, there is no centralized data store anymore; the business data is distributed to several places, as described above. Also, company data is often not stored exclusively inside the organization itself, but is distributed among subsidiaries or partner companies. Some data, such as financial data and address data, is often bought or leased from intermediate companies.

- **Government Data** - As governments are large organizations which have separate departments for each specific task, they produce data that reside in different places.

- **Scientific Data** - For large experiments, the data produced can no longer be processed locally (e.g. CERN [12], Folding@home [57]), and methods for distributing the data for remote calculation will be developed.

- **Data Finding** - Storing the data is only one aspect. Of course the data needs to be found and retrieved again later. While web search engines allow for more or less accurate searching of public data on the Internet, new trends are focusing on more personal data searching techniques, which function in the same way as existing web search engines (like Microsoft WinFS [50], Apple Spotlight [4] or Google Desktop Search [28]). There is also a trend towards a more structural storage of data for easier finding. Keywords are: semantic web [52], metadata [27, 82], relational data [86].

  - **Personal Data** - With the increasing amount of personal data (emails, multimedia data, documents) stored at multiple places (work, home, mobile devices), finding the necessary data is becoming more crucial and complex. Today finding information on the Internet is easier than finding data on a personal computer. As most users store the data in an unstructured way, the common methods for finding data - indexing the data before searching, as in web search engines - are normally used.

  - **Business Data** - Growing complexity and dependency on information technology is not easy to overcome. Companies normally use a more structural approach in managing data (in contrast to the personal data methods described above) and try to store data together with metadata or along with other information, even in a relational database. This allows for more exact searching capabilities.

  - **Government Data** - For several tasks, it is necessary to combine different data sources to get the desired result. For instance, if the immigration bureau needs to exchange its data with other countries, the data and metadata must be stored in a structured way. For exchanging those data with other authorities, standards for exactly how the data should appear are needed. The trend here is to have open standards for data exchange.

  - **Scientific Data** - The data itself is not always the interesting part. Oftentimes, data needs to be further analyzed and processed to get the desired

results. So finding the data means finding interesting relations between data. Here the data itself has to be very well-structured to allow automatic post-processing and also well-described so it is transparent and others can follow the experiments. Often for smaller projects, the data is not shared with other groups and it might be sufficient not to describe the data completely, but only to structure it for the needs of the project itself. For larger experiments, the data is shared with multiple groups and strict standards need to be put in place to allow other members to find the necessary data.

Generally speaking, if you need accurate access to data (e.g. finding *all* bank accounts within a specific criteria) the key is to structure the data in a precise way. If you are satisfied with inaccurate access to data, other heuristics can be used to give access to unstructured data (e.g. general indexing of plain text, as done by Google or others).

Quintessentially, it could be argued that:

"Storing data is easy - finding the data afterwards is the challenge."

### 2.1.2    Managing Data

For the storage and retrieval of data, there are several approaches, which can be distinguished into two main fields (as already described above):

1. **Structured Data** - This data has some sort of defined structure. This can be additional data describing this data (metadata [27, 82]) or other methods for structuring, such as grouping data and describing relations between data (relational model [86]). This structure can be added manually, e.g. by a person who has knowledge about how to structure a certain piece of data, or automatically, by using some heuristics to classify the data. Both methods have pros and cons, as a manual approach is always biased in favor of the person who does the manual classification and automatic categorization will suffer from inaccuracy because of limitations of the algorithms used. Examples are: XML files with a well-defined schema definition, relational databases and categorized articles.

2. **Unstructured Data** - This is data which has no structure. This can be data of any kind, which is mixed together with no relation between. Examples are: A storage medium with an arbitrary number of files of any kind, such as Word documents or pictures.

There are several approaches to managing this kind of data. Some of them allow for searching at different levels, while others only store the data and allow the user to browse all of it. Some examples (sorted from small to large systems) are:

- **Arrays and similar** - These are low-level, used by the application developer: These contain unstructured data and only the specific application can understand and use it. Even the computer memory (RAM) has some kind of array-like structure. The only way to access the data is by asking for the data at a

specific address. Example:
"Give me the data at position 5 of the array"

- **Key-Value Storage** - This is low-level, used by the application developer: The data can already be slightly structured: Each data item has a describing key, which then holds the corresponding value. Depending on the key-value relation, the data is more or less structured; e.g. if the key is only a random number, then it probably does not add any description to the data itself. In contrast, if the key contains a valuable description, such as "name" it adds some metadata to the data. Generally, there are two types of key-value storage mechanisms:

  - **Unsorted** - The keys are unsorted. This means the data is not ordered. This is mostly implemented as a "hash table". Example:
    "Give me the data for the key 'Zurich'"
  - **Sorted** - The keys are sorted according to a criterion (e.g. alphabetically, numerically). It is mostly implemented in a tree-like structure, with which the ordering is preserved at insertion time. Additionally to unsorted access, it allows queries such as:
    "Give me all data for the keys which start with 'Z'"

  For a more detailed description of key-value mechanism, see section 2.3.

- **File Systems** - From the low-level point of view, the data is structured in some way, e.g. the file system knows at which position on the disk a specific file is stored and in which sub-directory it is located, but mostly this information is useless for high-level usage; e.g. to find a file which contains a certain word, all the data on the file-system needs to be browsed. Example:
  "I have to browse the file system to find the file I need"

- **Databases** - Contain structured data. Furthermore, the database requires that only data that conforms to this structure can be inserted. Before filling a database with data, the designer has to decide how the data should be structured. Changing the structure afterwards is not an easy task, as adding a property mostly means updating all existing data items with these new properties. For accessing the data, there are sophisticated methods, which allow for fast and accurate searching of data. Example:
  "Give me the data which has the property 'city' set to 'Zurich' and which has the property 'income' greater than '100,000' Swiss francs"

- **Web Search Engines** - Depending on the architecture, these are based on databases or have their own techniques for storing the data. The data can be the web address of a page or even an archived copy of this web page [30]. To allow searching for unstructured data, the search engine maintains a huge index for the web pages to search. The algorithm to index and finally rank the index is unique to each search engine. Example:
  "Give me the link of a web page which is relevant for my word, 'Zurich'"

- **Peer-To-Peer Systems** - This is another unstructured way to store data: Several computers are connected to each other and each acts as a server and client at the same time and are able to retrieve some data. There are some large systems, such as the older Gnutella [24] network. Mostly those systems only allow for inaccurate, rudimentary and very slow searching capabilities (normally only on file names or some limited predefined properties such as file size or type of music encoding used, but all these additional properties do not have to be included all the time). Example:
  "Give me the data which has to do with 'Zurich'. This will not return all data, as some cannot be found, and can even return data which has nothing to do with 'Zurich'"

As already noted in the introduction, most systems described above are internally built on top of a key-value storage architecture and could therefore profit from a more general approach to do the distribution already on the key-value layer. Such systems are:

- **File Systems** - The base for most file systems is a sorted key-value storage architecture. Examples include: Microsoft NTFS [84] file system, Linux ReiserFS [54], Linux Ext2/3 [10], Apple HFS(+) [80] and others.

- **Databases** - Here also, key-value based solutions are commonly used as their underlying storage architecture, e.g. for building an index or allowing for sorted access to data [78].

- **Web Search Engines** - As some of them are based on databases or file systems, key-value kind storage algorithms are also found here, e.g. for the creation of the index [9].

- **Peer-To-Peer Systems** - Most of them are built on top of an unsorted distributed key-value architecture [24].

## 2.2   Distributed System

According to Tanenbaum [69]:

"A distributed system is a collection of independent computers that appears to its users as a single coherent system."

As it appears, this definition is quite difficult to implement, so normally the following weaker definition [44] is used:

"A distributed system is a collection of independent computers that are used jointly to perform a single task or to provide a single service."

Building a distributed system will always have some trade-off between different desirable properties. E.g. a system that is fully secure will probably not be fully efficient [43]. Also for specific use cases not all distribution problems need to be solved. Especially in our case, to build a distributed balanced tree, not all aspects need to be addressed. General important goals in a distributed system are (taken from different sources: [21, 22, 32, 43, 48, 60, 61, 62]):

- **Transparency** - For some applications it is necessary to hide from the user the fact that the system is distributed, so the user thinks he is working on a local system. There are different kinds of transparencies:

    - LOCATION - Users unaware of location of resources
    - MIGRATION - Resources can migrate without name change
    - REPLICATION - Users unaware of existence of multiple copies
    - FAILURE - Users unaware of the failure of individual components
    - CONCURRENCY - Users unaware of sharing resources with others
    - PARALLELISM - Users unaware of parallel execution of activities
    - OTHERS

- **Flexibility** - Depending on the need, a distributed system should be extensible, so additional components/services can be added/changed. Interface specification should be open, to allow easy addition of new features or enhancement of interoperability. Flexibility can also have an impact on other factors, such as reduced performance and security.

- **Reliability / Availability / Redundancy** - To enhance availability, load-sharing can be implemented with additional redundant data, so if part of the system fails, the system is still fully operable and after the failure, the whole system should stabilize again. There is also the problem that in a distributed system we have more components that could potentially fail, so to increase reliability, some recovery methods need to be implemented to bring the system back into a consistent[1] state after a partial failure. Reliability/Availability/Redundancy can have an impact on performance, flexibility, scalability, concurrency and basically all other elements.

- **Performance** - There are multiple measurements of performance, such as: Response time, throughput, system utilization and network capacity used. Any system will benefit from performance enhancements. In distributed systems, performance can have an impact on other desirable properties, such as transparency, flexibility, reliability and security.

---

[1]There are some objections that for a large distributed system, full consistency can't be reached and that, therefore, some ideas need to be implemented so the whole system will work, even if it is not in a fully consistent state (e.g. voting mechanism to check whether result can be trusted, multiple querying to see if the same result is retrieved).

- **Scalability** - "A system is said to be scalable if it can handle the addition of users and resources without suffering a noticeable loss of performance or increase in administrative complexity" [B. Clifford Neuman].
  There are three metrics in a scalable system. Depending on the needs, some of them are more important than others:

  - SIZE - Number of users and resources
  - GEOGRAPHY - Distance between users and resources
  - ADMINISTRATION - Number of organizations that exert administrative control over parts of the system

  Scalability can have an impact on performance. Sometimes a distributed system only provides better performance if a certain initial size of the system is being reached[2]. Scalability influences reliability, availability, redundancy and flexibility.

- **Concurrency / Parallelism** - Components in a distributed system are normally concurrent by design, as each component can perform its own calculation (those can be independent of or dependent on other components). Another aspect is the synchronization of information across the system. Concurrency has an impact on reliability, availability and redundancy as well as on performance.

- **Security** - A distributed system is more complex than a similar centralized one, so for components which need to exchange information, there is the risk that malicious components can harm the whole distributed system. Security has an impact on reliability, performance and scalability.

### 2.2.1   Distribution Methods

There are many methods to build a distributed system. Below, some strategies from large database vendors are discussed, how they distribute their data to multiple nodes. As those product portfolios are huge, and also because information is restricted, not all aspects are shown here. Currently the topic "distribution" is very hot, especially with regard to grid computing, and as vendors are trying to get into this area, the competition is quite tough.

Generally, a full-blown database (in contrast to our approach of a low-level database) allows several approaches to actually achieve distribution. E.g. the distribution could be done at table level, SQL level, or even higher at application level.

To compare the three databases, papers from [33, 37] helped give a general overview, although one of them is sponsored by IBM.

---

[2]Often a distributed system containing two nodes is slower than a comparable centralized system.

**2.2.1.1   Oracle - Oracle Corporation**   Oracle's solution for distributing data is based on using a "shared-disk" architecture [37]. In the most current version (Oracle 10g), Oracle is going into the direction of grid computing, but uses the common former distribution technique (RAC - Real Application Cluster) as the foundation of their grid offerings. RAC is a cluster database with a shared cache architecture that runs on multiple machines attached through a cluster interconnect and a shared storage subsystem [51].



Figure 4: Oracle 10g Grid Overview (taken from [56])

Using a cluster architecture has some limitations, as all machines need to be connected using a cluster-interconnect (a fast path from machine to machine); so distribution is limited to a high-performance LAN/Memory interconnect. Those machines can't be distributed over a larger area (e.g. the whole Internet). Also, depending on the cluster architecture used as a foundation, there are some more limitations. From our own experience in the telecommunications field where we deployed clustered systems, we have found that most architectures can only scale to a small number of nodes inside one single cluster (8-32 nodes per cluster [38]). The advantage of such a system is their high performance capability, because all machines are within a very fast network, but the limitation is the small size of such a system and the inability to distribute the cluster itself over large distances - so all nodes participating in a cluster must reside at the same location.

Using a storage subsystem (also known as shared storage, storage area network) is limited to a local distribution of the system, as access to this storage system has to be fast. The disks can't be spread over a large distance. Also, building a shared storage system can be quite complex, because it needs to have several mechanisms built in for concurrent access, e.g. a distributed locking manager, and therefore Oracle is limited

to systems, which actually support those features. See Appendix A on page 89 for an overview of two example architectures of such a storage subsystem.

The major benefit in using a shared disk/cluster architecture over a monolithic system is the increase in availability of such a solution. If one node goes down, all the storage is still accessible and another node can take over working on the same data.

Oracle has other offerings to distribute part of the database or the whole thing to larger areas. This is done in a replication manner using "Oracle Streams" to have multiple databases in sync over a bigger distance or using partitioning (explained below).

**2.2.1.2   DB2 - IBM Corporation**   In contrast to Oracle's offering, DB2 follows the "shared-nothing" approach: Each node owns and manages its own resources, including the data it has access to and the DB cache [37]. This approach was chosen to address the lack of good shared disk support (e.g. distributed locking mechanisms) in Linux, UNIX and Windows. Using this approach, IBM is claiming to reach "infinite" scalability to at least 1000 nodes [33]. Generally, this technique is closer to our b-tree distribution method, where we also plan to have "independent" nodes.

To distribute data in a DB2 environment, it mainly uses partitioning (MDC - Multi Dimensional Clustering [34]) to spread a table across multiple partitions, residing on multiple disks or systems (so a relational database table will be split horizontally into several chunks and those are then placed into different nodes).

Figure 5: Table Partitioning

This also allows for parallel executing of queries, as every node can process its own data and only the result will then be combined. Generally this technique is also available in Oracle.

**2.2.1.3   SQLServer - Microsoft Corporation**   In contrast to Oracle and DB2, SQLServer does not support clustering of machines by itself. It is limited to one machine, although it can be a multi CPU machine. Distribution of data is only possible using "federated databases" (linking distributed, heterogeneous databases together using software), [33] e.g. by having a database installed in each country and then running an application which will query each database independently and combine the results. This approach is similar to partitioning, but in the case of SQLServer, the user has to do it by himself in the application accessing the database. This approach has several limitations. E.g. in case of renaming parts of the table, this renaming has to be propagated manually to all participating database nodes.

## 2.3   Local Data

Regardless, for monolithic or distributed systems, they need to store data. As each machine only has direct access to its own hardware, the actual low-level storing is always done locally[3]. This means that the actual algorithm to store the data can be the same for monolithic machines or nodes of a distributed system.

There are several algorithms for storing the key-value pairs. In this section, we will mainly focus on algorithms for sorted storage of key-values and will further focus on the most commonly used algorithms.

### 2.3.1   Binary Search Tree

A binary search tree [76] is a special form of a binary tree [77], in which all nodes are inserted in an ordered manner: If you pick one node, all sub-nodes to the left always contain values less than this node, and all sub-nodes to the right contain values greater than this node itself, as shown in Figure 6. This data structure is the base for most others, which just add functionality to it.

---

[3]The terms distributed/monolithic have no exact boundary, as one could argue that a single machine itself has distributed components: CPU accesses the memory over an internal bus-system and accesses the disk over an external cable, for example. Nevertheless, in this work distributed means connected using a network interconnect to each other. All other bus-systems or internal communication paths are regarded as part of a monolithic system.

Figure 6: Binary Search Tree

As a binary search tree does not rebalance itself during insertion of new elements, the tree might, in the end, not look like a tree but more like a linked list, where all nodes are chained after each other. In this case, the tree is highly unbalanced.

Figure 7: Highly Unbalanced Binary Search Tree

This will limit the search performance of such a tree, as for the worst case scenario described above, the search-algorithm basically has to look up all nodes of this tree and is of order $\mathcal{O}(n^2)$, where $n$ is the number of nodes of the tree.

### 2.3.2   Self-Balancing Binary Search Tree

A self-balancing binary search tree [87] is a special form of a binary search tree in which the algorithm tries to limit the shortcoming of a normal binary search tree by implementing a method to prevent the unbalanced growth of the tree. The height of the tree is automatically kept to an optimal minimum. Depending on the actual implementation, the worst-case lookup time can then be reduced to an order of $\mathcal{O}(\log n)$, where again $n$ is the number of nodes of the tree.

Some data structures that use a self-balancing binary search tree are:

- **AVL Tree** - Invented by Adelson-Velskii and Landis [73], 1962.

- **Red-Black Tree** - Invented by Rudolf Bayer [85], 1972.

- **Splay Tree** - Invented by Daniel Sleator and Robert Tarjan [88], 1985.

- **B-Tree** - Invented by Rudolf Bayer [74], 1972. This is the most commonly used form of a self-balanced binary search tree. There are also a variety of slightly different algorithms which all are based on the original b-tree.



Figure 8: B-Tree

As shown in Figure 8, all leaf-nodes have the same height of 3 and always need the same amount of time for lookup. In contrast to a normal binary search tree as shown in Figure 6, the height depends on the nodes: Nodes 7 and 24 have a height of 3, while the other leaf-nodes have a height of 4. Therefore, the lookup for each leaf-node takes a different amount of time. Hence, the tree could be even worse - a few leaf-nodes could have a much greater height, for instance (as seen in Figure 7).

### 2.3.3  Skip List

This data structure is also commonly used and is based on parallel linked lists and was invented by William Pugh [58], 1990. Instead of only linking the successors/ancestors, crosslinks from one element to another arbitrary one are possible.

Figure 9: Skip List

It has advantages in parallel computing where the re-balancing (or re-linking) can be done independently of the other nodes: Each element can decide by itself to whom a link should be established, without global communication. Nevertheless, it has some shortcomings, as there is still a possibility (although very low) that at the end the skip-list is just one huge linked list, which would then have a worst-case lookup order of $\mathcal{O}(n^2)$. But in normal circumstances, this algorithm can achieve a lookup performance similar to a balanced tree of $\mathcal{O}(\log n)$.

## 2.4   Distributed Data

While good implementation of balanced trees for centralized architectures exist, we are looking for a solution on how to distribute sorted key-value like structures (e.g. balanced trees) among multiple nodes, to gain higher throughput, as well as redundancy. With this approach some new problems, which are normally not found on a centralized system arise:

- **Data Balancing** - Data is now located on different computer nodes, so a mechanism that will distribute the whole database over all participating nodes has to be built. Depending on the needs, some nodes can hold more data than others (e.g. if they have larger disks, more processing power) or data should be evenly distributed among all nodes.

- **Data Reallocating** - A database is highly dynamic, so it has to move its data automatically among other nodes if needed. (E.g. a lot of inserts could lead to a growth in some limited number of nodes, so data has to be transferred to other nodes to gain good balancing data). This is related to data balancing.

- **Communication Between Nodes** - Depending on the design decision, internal communication between participating nodes will be needed (e.g. to reallocate data). Those have to be reduced to a minimum, or at least be on local order only. E.g. there is only communication between "neighboring" nodes, not broadcasting to all nodes.

- **Access Time Concerns** - In a distributed environment, all accesses (e.g. searches, inserts, deletes) eventually go over a network connection, depending on architecture: An access-request could even be delegated from one node to another, which could have an impact on access latency and throughput.

### 2.4.1   Small-Sized Systems

This section describes some methods for smaller systems of how the data could be distributed and still allow for sorted key-value access.

**2.4.1.1   Paging B-Trees**   One promising approach is the use of paging b-trees [55], where sub-trees are relocated to different nodes. To prevent splitting the tree into too many small sub-trees, the distribution is done on the depth-level of a tree. E.g. if the whole tree has a depth of 4, then splitting is done in such a way that 2 depth fields are always grouped together. This also solves the problem of locality, because now local/similar data is most probably on the same node and not spread to other nodes as well. Experiments have shown that the overhead in querying these paging b-trees is around 5%, while the overhead for inserting is around 50%, which seems quite high. It seems it is better suited for applications, which read more from the tree than write to it.



Figure 10: Example of a b-tree split into several parts and put on multiple nodes (taken from [55])

**2.4.1.2   Concurrent B-Link Trees**   A more theoretical work on distributing b-trees is done by Johnson/Colbrook [39]. As a foundation, a concurrent b-link tree [47, 63], which, according to the authors, has the highest performance of all available b-tree algorithms, is used. This is done with fine granular locking mechanism in case of concurrent access on the b-link tree. Further b-link trees help in finding the

data, even when parts of the upper tree are missing, because all leaves are linked to each other, so in case of a node failure or a short-local inconsistency (e.g. during a balanced tree re-organization process) the data can still be found.

Because all b-tree queries normally start at the root node, Johnson/Colbrook [39] proposed to replicate the root to all available nodes to have better throughput. The same is done for sub-trees, so in the best case the request will never leave the first chosen node and will find the data without hopping to different machines.



Figure 11: Shows the replication of sub-trees to get higher throughput (taken from [39]). The numbers show, on which CPU this part of the tree is stored. E.g. the root tree is stored on CPU 1, 2, 3 and 4.

### 2.4.2   Medium-Sized Systems

The difference between small and medium sized systems is not very big. Most people consider a system medium sized if it is too large for huge data exchanges (e.g. it does not use a shared disk approach) but it is still efficient enough to broadcast messages to all members of a system. Mostly those systems are also located at the same site and are not distributed over the whole world.

**2.4.2.1   Boxwood Project**    This project is part of a Microsoft research program [49] to build a distributed storage architecture. It focuses on building a distributed file system on top of a distributed data structure. To solve this problem, they first want to build a distributed balanced tree architecture and then use this as a foundation for their file system. Unfortunately, no source code or architectural overview, on how they actually want to build it, has been released so far.

### 2.4.3   Large-Sized Systems

Large systems are considered Internet-scale systems, which also have the characteristics of unpredictable failure of nodes and communication between them. Such systems are mostly built around a structure common to P2P (Peer-to-Peer) networks or have a hierarchical approach (e.g. such as DNS systems).

**2.4.3.1   Structured Peer-to-Peer Overlays**   Structured Peer-to-Peer (P2P) over-
lays are the foundation of highly distributed P2P networks and currently a hot topic.
In contrast to older P2P networks like [24], those new approaches [23, 40, 59, 66, 90]
focus on true decentralized mechanisms, scalability, self organization, higher reliabil-
ity and better performance. They allow routing to location-independent names (e.g.
"route message to *key*, where *key* is located on a node-not-yet-known in the network).
On top of such an architecture, higher-level operations, e.g. a distributed hash table
(DHT), decentralized object location and routing (DOLR), group anycast/multicast
(CAST) and possibly, also a distributed balanced tree (DBT) can be built.



Figure 12: P2P overlay (key-based routing layer) with several services on top of it
and applications using those services (taken from [17])

There are some projects going on which are working on a distributed low-level
data-structure called Distributed Hash Tables (DHT), which are commonly used in
Peer-to-Peer network architectures and are also found in some research projects re-
lated to distributed storage. The difference between hash tables and b-trees are that
hash tables only allow unsorted access to keys (e.g. a lookup of a range of values is
not possible) and b-trees have sorted keys. It seems that for distributed storage (e.g.
file-based access, Peer-to-Peer networks) a DHT is enough, as several projects [8, 42]
are using it. Nevertheless, no implementation of a sorted distributed data structure
has been found so far.

All those P2P overlay projects focus on large-scale distribution, and therefore try
to minimize global communication, so we can learn from them:

- How those projects do data balancing.

- There are some techniques to build a b-tree on top of a P2P overlay (e.g. [13]).
  With this approach we could use an already proven overlay implementation
  and add a b-tree to it. It has to be shown whether there are drawbacks. E.g.
  data balancing could be worse or access time could be less than perfect.

- There exist already a DHT implementation on top of those P2P overlays. Per-
  haps we could find a way to build a distributed b-tree on top of such a DHT.

**2.4.3.2   Sorted Distributed Hash Tables**   The main advantage of using a balanced tree instead of a hash table is the fact that all keys are sorted. Besides using a b-tree, there are other techniques, e.g. those based on hash tables, to get access to sorted keys (some are described at [89]). With those methods, we could use existing DHT implementations and build a system, which has the same functionality as a distributed balanced tree.

One proposed data structure, which would provide the full functionality of a distributed balanced tree are skip graphs [7]. Skip graphs are based on skip lists [58] (an alternative to balanced trees, see section 2.3.3) and enhance them with DHT to get a truly distributed balanced tree. Nevertheless, Yalagandula [89] showed that skip graphs do not seem to scale well.

Another is Squid [64], who basically allows multidimensional sorted keys, but here also Yalagandula [89] showed that this method does not scale well.

Yalagandula then proposes his own approach by combining Tries [1, 25] with hash tables by first filling the data into a Trie and then mapping this structure onto physical nodes using DHT techniques. Although the author plans to implement a prototype of his idea, no such implementation has been found so far.

**2.4.3.3   P-Tree (P2P Tree)**   Instead of splitting a single balanced tree across all nodes, there are also ideas to have multiple independent b-trees and assign them to each available node. One such an approach is described in [13], where they introduce P-trees (Peer-to-Peer trees) and some kind of semi-independent B+-trees (a flavor of the common b-tree). This is a mixture of having one full tree distributed over all nodes and having several independent trees.

Each node will build up some part of a global index where the actual data can be found. This index does not need to be fully consistent; in case there are some failures, the search will take just a bit longer. The main disadvantage of this structure is that the data, which can be stored on one node, is very limited (e.g. only one single item of data can be stored - this limitation can only be overcome by starting multiple virtual nodes on each real node, but even with this way, the amount of data to be stored is very limited). This data structure can be seen more as an indexing framework, where only a very limited amount of data needs to be stored at each node. Or the node itself is the interesting data to be found, so this data structure could be used for node discovery, e.g. as a "replacement" for DNS. (It would then be called DDNS, Distributed DNS).

**2.4.3.4   P-Ring (P2P Ring)**   To overcome the limitations in P-trees, an extension has been proposed: P-ring[14]. With this approach, it is now possible to store large amounts of data on each node and the current limitation only seems to be that the amount of data stored on each node is expected to be roughly the same. But also here, no prototype has been released to the public so far.

## 2.5  Performance Aspects

When designing a distributed system, special care needs to be taken with performance issues. Generally, a distributed system will send some information between each of its nodes. Depending on the actual architecture, such messages could be:

- **State Information** - General alive messages, state-change information or messages needed for synchronization of the overall system wide state.

- **Managing Information** - Such as information about joining, leaving, or statistics.

- **User Information** - The actual user messages, the data packets that need to be processed by the system.

- **Overhead Information** - As each network layer, including the application protocol, adds some housekeeping information to each message, the overall message grows in size.

Sending messages and accessing data remotely is always very expensive, as in current network topologies (e.g. Ethernet, IP) transmitting information between nodes is slow, compared to local access methods, as shown in Figure 13. This is especially true if the nodes are distributed on the Internet, where the network latency is quite high compared to the other methods.



Figure 13: Access Time Differences (taken from [53])

### 2.5.1  Minimize Disk / Network Traffic

The goal is to minimize disk- and network-related traffic, as those have the biggest impact on slowing down a distributed system. As disk access is always done locally, it is not only a problem in distributed systems but also in basic monolithic ones, and several implementations for local key-value storage have already addressed this issue (e.g. the BERKELEYDB [65] implementation of a balanced tree implements advanced caching and disk access procedures to further reduce disk access). For discussing performance issues on a distributed system, minimizing and optimizing network traffic becomes the main issue. Basically this problem is split into two main parts:

- **Problems you can control** - These are issues which can be reduced by choosing a good design and implementation of the distributed systems. Such things are:

  - **Algorithms** - Using efficient algorithms probably has the highest effect on reducing network traffic - the goal is to reduce all unnecessary traffic and also limit the amount of data to be exchanged. This has to be done not only for the network-related part, but also for the underlying architecture. E.g. if the system can also correctly operate in a not fully consistent state, the number of messages related to consistency can be reduced. Or by taking advantage of parallelized algorithms, where the actual workload can be split into multiple subtasks, which could then be processed independently, the latency of each request can be reduced.

  - **Broadcasts** - Instead of sending requests to all members of a system, the message could be sent only to a subset of members (multicast) or the number of such messages could be reduced in favor of more personalized messages (point to point).

  - **Caching** - Requests could be cached, so if a similar access needs to be processed, it could be found in the cache and would not need to be sent out to the network.

  - **Overhead** - Each network layer adds some administrative overhead to a message - generally this is out of your control, but by wisely choosing the right network medium (e.g. TCP or UDP) and the serialization coding (to convert internal data to binary data) and not adding much overhead with your own application, this overhead can be reduced.

- **Problems which are out of your control** - These are issues which arise in the network and are not in control of the application itself. These issues harm the overall system performance, as the system needs to react to such failures with additional administrative effort (e.g. by implementing fault tolerance mechanisms). Some characteristics of problems are (taken from [11]):

  - **Failures** - Network components can fail at any time. This is especially true for larger systems, where the probability that some parts of the networks are down is higher than on smaller systems. Besides network failures, nodes participating in the distributed system can fail for any reason.

  - **Latency** - The time a network message needs to pass from one node to another can vary, and this is especially true for larger systems, or if the transport channel is also used for other traffic (e.g. a company wide intranet, which transports emails, web page requests and also the communication messages for the distributed system).

  - **Throughput** - The amount of data which can be sent during a specific time will also vary, especially for larger systems where the network paths have different characteristics with regard to speed.

### 2.5.2   Asynchronous Processing

A general problem for systems which access different components which all have different access time constraints (as shown in Figure 13) is that faster components could process much more data but are slowed down because they need to wait for slower components. E.g. the disk or network part of processing will slow down the CPU because it has to wait until the data is stored to disk or sent over the network and waits for a response. The CPU could do other work at the same time instead of waiting for responses.

There are a few methods for overcoming these problems [72]. What they all have in common is that they don't wait for a full request to be processed. They accept new requests and try to execute them while waiting for the slower components to be finished.

**2.5.2.1   A Thread per Request**   This method basically assigns a thread for each incoming request, to process it. As the number of requests increase, the number of running threads increases as well, as shown in Figure 14. Generally, each thread processes the whole request from start to finish by itself and then sends back a response to that request. Similar to old manufactures, where the employee is building a product from start to finish.



Figure 14: A Thread per Request

This method allows for the processing of multiple requests at the same time. Even when one request takes a long time to process, another thread can be processing another request. This method has the drawback that a lot of threads are running at the same time and the operating system needs to switch thread context quite often. As

the number of threads increases, the administrative overhead to manage those threads also increases and overall performance suffers.

There are variations to this approach, such as limiting the number of threads and only allowing a few threads to run at the same time. Additional work, which cannot be processed because no idle thread is available, will be put into a queue and every time a thread has finished its work, it will retrieve the next pending request from its queue. This allows for more control over the number of parallel threads running on the system.

Using a multi-threaded approach, great care must be taken in synchronizing shared resources. If different threads want to access the same resource (e.g. file system) locking must be put in place. This leads to lower performance, as this synchronization/locking overhead needs more CPU resources. In addition, the complexity of the program can increase, as multiple threads could access each component at the same time.

**2.5.2.2   Event-Driven Architecture**   Instead of letting a thread process the whole request from start to finish, event-driven architecture basically runs around events. Events are more finely granulated than a whole request. E.g. a request is a full-blown application-related task, such as: "Fetch web page with name seda.html", while an event is just: "I have a TCP-packet on my network interface", so the whole request is split into multiple events, and for each event type there is a responsible component, which processes it and also generates new events which are then processed by other components. The communication between those components is separated by event-queues; this allows for nicely decoupled components. Similar to more modern assembly lines, where each employee is only building a part of the whole product - this is normally much faster than old manufactures.

The number of threads is not controlled by the components themselves, but by a scheduler component, which gives the thread of execution to those components that have full event-queues.

Figure 15: Event-Driven Architecture (taken from [72])

As Figure 15 shows, an incoming request will trigger one or more events, which are then passed to the responsible event-queues for processing (e.g. the request first needs to be parsed). This component itself will trigger new events for further processing (e.g. the parser will now check whether the request could already be answered with information from a cache).

A possible implementation of such an event-driven architecture is SEDA [72] (Staged Event-Driven Architecture), which is used in several applications already. Beside providing access to an event-driven architecture, it also provides additional useful functions for an application programmer and implements a sophisticated thread pooling mechanism which allows for controlling the number of threads for each component (called stage in SEDA) depending on system load.

As we will see later, the proposed prototype ERT also uses this kind of architecture and borrows some ideas from the SEDA approach (e.g. the components only communicate through queues and never call methods directly).

# 3   Overview of ERT (Extended Root Tree)

When sorted access to data is needed in a centralized environment, b-trees (balanced trees) and their similar flavors are commonly used. We will now discuss how such a data structure can be scaled to a distributed environment to provide the same functionality as a centralized b-tree. The method discussed here extends the algorithms found on balanced trees or other local sorted key-value algorithms with additional functionality to get the desired distribution. One can say that the data structure in the end is not a distributed b-tree, but solves the same problems and gives the same interface to key-value data:

- It allows access for exact match results (e.g. lookup of a single key).

- It allows range queries (e.g. lookup of a range of keys).

- It implements *put(key, value)*, *get(key)*, *getNext(startKey, amount)*, *getPrev(startKey, amount)* and similar interfaces to set and get data from the distributed data structure.

The proposed high-level architecture is designed to scale from small systems to very large systems. The main idea behind this is to reuse as much as possible already available technologies and software, to lower the risk of having to build a completely new system. This also has the advantage that if new trends or new problem-solving techniques come up, it should be possible to incorporate them into this system. Another point is that for different requirements (e.g. smaller systems, larger systems) some components can be tuned for their specific environment, without changing other parts. Such a framework is shown here (similar frameworks can be found in [15, 17]):



key transformation      node lookup      data store      replication

Figure 16: The different components for building a distributed data structure.

The next part of this article shows an overview of how everything works together, and the following part goes into a bit more detail about each component.

## 3.1   Basic Concepts

The main benefit of a distributed b-tree over a distributed hash table is the sorted access to its data. One idea would be to store the keys out-of-order (unstructured) and only to restore the order during query-processing time. But because the aspect of sorting is the fundamental feature of the whole system, it makes sense to sort the keys at insertion time, so at all times, the whole data structure is in a sorted state (structured).

Considering we have a user key-space $[A, Z]$, this range has to be partitioned and distributed among all participating nodes[4]. Distribution can be done at key-level, so each single key is managed on its own (e.g. [13]), or we can group some keys together (*"key-group"*) and then manage the whole group as one single entity. A *key-group* is like a very small range of the whole user key-space and all keys inside such a *key-group* $a \subset [A, Z]$ are ordered and are responsible for exactly this assigned range. No other *key-group* $b \subset [A, Z]$ contains a key which is also contained in *key-group* $a$ (so $a \cap b \cap \ldots = \emptyset$). Each *key-group* has a unique label (id, hash-key) to identify it[5]. Furthermore, a range of continuous *key-groups* will be called *key-partition* e.g.: $\alpha = [a, d]$. As an overview of how the mapping of a user key to an internal key (= *key-group*) and the mapping to a node is done, take a look at Figure 17. Figure 3 will show, how such a structure is distributed among multiple nodes.



Figure 17: Mapping of user key-space $\rightarrow$ internal key-space $\rightarrow$ nodes.

---

[4]At this point, we won't look at data replication. So if one node goes away, all data for which this node was responsible are no longer accessible.

[5]In this paper, elements of the user key-space are written in capital letters (e.g. A, B, C) and *key-groups* are labeled with lowercase letters (e.g. a, b, c) and *key-partitions* are labeled with Greek letters (e.g. $\alpha, \beta, \gamma$).

Figure 18: Transformation of keys to our internal key-space.



Figure 19: Range $[a, o]$ is distributed among 4 nodes. Node 1, 2, 3 hold two *key-partitions*, node 4 only one.



Figure 20: All sub ranges are linked to their successors and predecessors.



Figure 21: User keys C, P and Z are stored on the corresponding nodes.

As the user of such a distributed data structure wants to have full flexibility with the keys, he might use alphanumeric, large range or small range, but the keys can't

be used as received from the user. They must be transformed to an internal key-range (example in Figure 18).

Assuming we have a whole internal key-range (= *key-partition*) of $\alpha = [a, o]$, then each participating node will be responsible for one (or more) non-overlapping sub ranges of the whole key-range (see Figure 19: E.g. node 1 consists of sub ranges $\alpha = [a, c]$ and $\vartheta = [j, k]$) and all nodes will know where the successor and predecessor of their sub ranges is (Figure 20).

The transformation process depends on the size of the system. For example, on a small 4-node system, it would be possible to have 4 large *key-groups* and just establish a 1:1 relationship between *key-groups* and nodes, so each node is solely responsible for one single *key-group*. On a large system, there are many more *key-groups* and each contains only a very small subset of the whole user key-range.

To retain the original user key, it is stored on the node itself (Figure 21).

The commonly used operations are described here:

### 3.1.1   Lookup / Storing of Key-Value Data

1. The user key is transformed to the internal key-space using a linear hash-function.

2. The corresponding node for this internal key is now searched using already proven P2P lookup methods (P2P overlays, distributed hash tables [17, 23, 40, 41, 66, 90]) or for smaller systems a simple lookup-table (static table, Zeroconf implementation à la "Rendezvous" [2, 3, 35] or like a DNS system).

3. The corresponding node will receive the request and will return/store the data.

### 3.1.2   Lookup / Browsing a Range

1. The lower end of the user key is transformed to the internal key-space using a linear hash-function.

2. The corresponding node is now looked up.

3. The node returns all values in the requested range.

4. The node forwards the requests he cant fulfill to the next node (successor) and so on.

For optimization, range queries could also start at both ends, so the lower-end and the higher-end are executed in parallel. This will require some kind of transaction ID, so if the two split and processed meet at one node (e.g. in the middle of the range), this node will know that it has already processed this query and will stop.

### 3.1.3   Node Insertion

1. The new node will randomly send out some requests to get the load of some nodes.

2. The new node will then contact the node with the highest load to migrate some data.

3. The new node will then register its successor and predecessor. (This can be done in the way b-link trees register new tree nodes [39]).

4. The new node will then officially be the holder of the *key-group*.

### 3.1.4   Node Failure

1. If a node leaves the system (planned or unplanned), the links to its successor and predecessor will be broken.

2. Its successor and predecessor will search each other and will link together. The range of the lost node will no longer be available (unless there was a data replication mechanism in place), but the rest of the distributed data structure will still be functional.

## 3.2   Key Transformation

The main idea behind key transformation is that our implementation will have an already pre-defined key-range for the internal key-space. Depending on the technology used for the node-lookup, this can be a very large integer space. If we used "Bamboo" [59] as the underlying node-lookup infrastructure, the internal key-space could be integers from $0$ to $2^{160}$. At the same time, the user key-space could be phone numbers ranging from 1,000,000,000 to 9,999,999,999 and therefore would need to be mapped to our internal key-space. Note that not the whole internal key-space needs to be used. In this example, the phone number range could be mapped 1:1 into our internal key-space or, better yet, to an even smaller internal key-space, because otherwise each *key-group* would only contain a single entry[6].

   This transformation needs to keep the ordering of the user key-space, which is why a linear hashing function needs to be chosen. The hashing function does not need to provide perfect ordering. If the hashing function only provides a proximity result (e.g. it finds that the requested key should be in *key-group* $a$ but is in fact in $b$), then the corresponding node, which holds the *key-group,* will forward the request to its successor node to query there. Of course, this will have an impact on performance, so the better the hashing function, the better the lookup performance[7].

---

[6]E.g. the range 1,000,000,000 to 9,999,999,999 could be mapped to 0 to 1,000. This way each *key-group* would hold up to 9,000,000 numbers and we could still scale up to 1,000 nodes.

[7]Forwarding data to the successor node only needs to be done for range-queries. For exact queries, the hashing function will never be inaccurate, because when the hashing function decides in which *key-group* to store the value, it will be the same as when the user queries for the same value for a lookup.

## 3.3    Node Lookup

Once we get the internal key for the request, we have to look up the corresponding host, which contains the requested key. All *key-groups* are distributed among all participating nodes. As a general rule, if we have several *key-groups* located at the same host, it is preferable that those *key-groups* form a closed range of the internal key-space (this is called *key-partition*) (e.g. it is preferred to have *key-groups* $a, b, c, d$ on one node and not $a, k, r, z$. This will help prevent unnecessary jumps from one node to the other during range queries).

As mentioned above, depending on the size of the system, different approaches can be taken:

### 3.3.1    Small-Sized Systems

For smaller systems, a static routing table can be put in place, e.g.:

| Internal Key | IP Address |
|:---:|:---:|
| a | 192.168.0.10 |
| b | 192.168.0.5 |
| c | 192.168.0.9 |
| d | 192.168.0.5 |

Figure 22: Static Routing Table

### 3.3.2    Medium-Sized Systems

In a more dynamic environment, in which nodes are constantly added/removed, DNS-like approaches or other local-network methods like Zeroconf (also known as "Rendezvous") [2, 3, 35] can be used. With these methods, every node broadcasts its internal key on the network, so others can easily find it. Zeroconf claims that their broadcast mechanism is resource limited and will not overload the network with unnecessary broadcast messages. Nevertheless, this mechanism only works on the IP broadcast scope and is therefore limited to the local network, as broadcast messages are not forwarded by a router to a foreign network.

### 3.3.3    Large-Sized Systems

Obviously for a large-scale systems, other mechanisms need to be considered, even non-centralized ones (e.g. for current Peer-to-Peer systems such as Gnutella [24]). In such systems, lookup becomes the most critical part of the system, as it is normally the slowest. Recently, there have been many projects working on a distributed lookup

architecture (sometimes also called "Peer-to-Peer overlays" or "P2P overlays" for short) [15, 17, 23, 40, 41, 59, 66, 90]. Some of them could be used for our approach.

Depending on the underlying system we choose, additional optimizations can be implemented. E.g. sometimes the key-space is not uniformly distributed - there are some spots ("data-spots") with a lot of data and some with no data at all ("data-holes"). In such a sparsely occupied key-range, not all *key-groups* need to be pre-allocated. They could be allocated as soon as some data fell into this *key-group*. This way, during an insertion of a new key, the system asks the nearest node to also add this new *key-group*. (E.g. if you already have *key-groups* $a$ and $b$ on node 1, now a new key, which would fall into *key-group* $c$, should be added. The system will ask node 1 to also hold *key-group* $c$.) This approach has two advantages:

1. The *key-groups* do not need to be pre-allocated, so use of the nodes is made more efficient.

2. Allows better scalability, as newly added nodes can efficiently be filled with new keys without having to migrate data from existing nodes.

To accomplish this feature, the underlying structure needs to support the routing of lookups for non-existent *key-groups* to the nearest existing *key-group*. (E.g. if *key-group* $a$ and $f$ exist, and then a lookup for *key-group* $c$ is done, the underlying structure needs to route this request to either $a$ or $f$. This only has to be done for range-queries. For exact queries, the system can immediately return, as the key does not exist).

## 3.4   Data Store

Each node has its own storage where all the keys of these *key-groups* are stored (see section 2.3). For efficiency reasons, each node only needs to have one single storage, and all key-ranges can be stored in the same place. As this storage also needs support for sorted access to keys, conventional b-tree (e.g. BERKELEYDB [65]) or highly concurrent b-link trees [39, 47, 63] can be used.

To make the system as flexible as possible and to allow further extensions, each key-value pair can have one or more pieces of additional metadata information, which is stored alongside the data. Such metadata information could be:

- Permission flags, so access to this data can be restricted to some users

- Time and date, transaction ID of last change

- How many times the data item has been accessed already

- Locks for transactions

- A flag if actual data is locally stored or stored on another machine of the network

This metadata is also used by the system itself for monitoring, development and tuning purposes. Beside having metadata and counters on data level, also other counters should be put in place on node level, e.g. the amount of network traffic coming in or out, how many queries the node already has answered, etc. This allows for monitoring of the whole system during testing, finding the best way to implement all its features. Later it can also be used for self-maintenance, e.g. to detect an overloaded system or cache frequently used data.

Each request coming to the node should have a "unique" transaction ID, so if the whole query is split in parallel and would reach a node twice - for whatever reason - the node will know that he already processed this query and does not need to do any further action. Those transaction IDs do not need to be stored for a long time, and can be removed after some time.

## 3.5    Replication

So far, we have not discussed replication of data in such a system, as there are several aspects and granularity to look at. E.g. replication could be used as an active/standby mechanism, so the replication is only used if the master is not available anymore, or in an active/active manner, where writing/reading could be done on either master or replica.

### 3.5.1    Active / Standby

This is the simplest approach and would only give the advantage of having a backup copy in case the master is not available anymore. Additionally, several replicas could be used to achieve even higher redundancy in case of failure. This method would not give any performance advantage of having multiple data sources for the same data, so access to its data could not be load-shared between master and replicas.

But the advantage is that this approach would not need to block write requests until all replicas could be updated, so this is rather simple to implement.

To include such mechanisms into our distributed data structure, a simple approach can be used: The original *key-group* acts as master and in case of an insertion/update of a key-value it propagates these changes to its first replica, and this first replica then propagates the changes along to the second replica (e.g. $a \rightarrow a' \rightarrow a''$). Alternatively, the master could also update all its replicas (e.g. $a \rightarrow a'$ and $a \rightarrow a''$). This could be done synchronously (the master waits to send back the response until all replicas have been updated, and then returns its status) or asynchronously (the master immediately returns its status, and updates its replicas later).

The master and its replica always send some tokens to probe their availability[8]. Alternatively, replica 2 could also only probe replica 1, in case of a hierarchical setup.

If the master is unavailable, the replicas will fail to send its probe tokens and the next available will be elected as master. If we have an underlying node-lookup

---

[8]Only the replicas have to send such a token. The master itself never needs to send a token to its replica, as he knows their status anyway, when an update fails.

Figure 23: Replication: A master and two replicas.

architecture as described in section 3.3.3, the node-lookup will automatically find the correct replica, as it cannot find the original host and will just forward the query to the next available *key-group*. (E.g. if a lookup for *key-group* $a$ fails because node 1 is not reachable, the underlying node-lookup architecture tries to find the next available *key-group*. In Figure 23 this would be *key-group* $d$ located at node 2, where the first replica $a'$ is also located.)

If the old master is available again, it will contact the current active replica to update its value and will then act again as normal master.

### 3.5.2    Active / Active

With Active/Standby, we gain redundancy in case of a node failure. With Active/Active, we could additionally gain load-sharing, although this method is more difficult to implement.

If some data that fall into the same *key-group* is "hot" and needs to be accessed from several clients at the same time, a potential bottleneck will occur, as only one single node is responsible for this data. Here replication could also help, as master and replicas would allow access to their data at the same time. Using this method, synchronization and locking between master and replica must be put in place, so all update requests must be atomic, and during a write/update, no client is allowed to access the same data for read/write. Depending on the transactional need, there are possibilities to tune the locking behavior, but nevertheless there is always some amount of administrative cost to update all replicas.

Generally we can assume that writing data will be slower, as all replicas need to be updated and during that time no other client is allowed to access the same data. Each time before updating the data, the master has to set locks on each replica and can then update the values.

Reading from those data can then be done in a load-balanced manner. The node-lookup mechanism will randomly (or any other load-balance mechanism can be used) forward the queries to either the master or his replicas.

Instead of using a separate key-space for all replicas, the node-lookup will just add a fixed prefix to its internal *key-group*. This prefix has to be large, so the chance of putting the replica on the same node as the master is minimized. Assume a prefix of 1000 and look at the master *key-group* with ID "35". Doing a lookup of "35" will

return node 1. As *key-group* "36" is probably also on node 1, we now take our prefix
of 1000 and then look up *key-group* "1035". This one is most probably not on node
1, and on this node we could store our first replica. The second will be on the node
with *key-group* "2035", and so on.

# 4   Implementation of ERT (Extended Root Tree)

This section will describe the low-level and implementation details for the distributed data structure. It will introduce the main concepts of the application and explain why a specific approach has been taken. There are two main roles a user of the system can have:

1. **End User** - This is the user of the Application Programming Interface (API) of the system.

2. **Administrator** - This user configures the system and does administrative tasks (like adding new nodes and removing nodes).

It should be noted that in some circumstances, both roles could be executed by the same person.

According to section 2.2, the system will try to fulfill the requirements for a distributed system. The following list summarizes the goals achieved or those still open for the implementation of ERT:

- **Transparency** - The current design tries to be as transparent as possible to the end user. It is important to know that the administrator of the system knows that the system is distributed, as he needs to configure the nodes. On the other hand, the end user API (Application Programming Interface) abstracts the fact that the system is a distributed one. There are different kinds of transparencies:

   - LOCATION - The end user is unaware of the location of a resource; he just has an interface to get/put/delete values and does not know where the actual data is placed (local/remote).

   - MIGRATION - The end user is unaware if resources have been migrated or are migrating at the time the user sends requests. In addition, migration does not alter the resource name or anything else.

   - REPLICATION - The end user is unaware that data exists as multiple copies distributed among some nodes. Also, the end user has no influence to stop/start migration for a particular data item. Replication can be

turned on/off by the administrator of the system. In the current version, replication is not implemented, but the design has already taken replication into account.

- FAILURE - Depending on whether replication is enabled/disabled, the end user will not see any disruption in service. The failure rate also depends on the degree of replication and the overall size of the system.

- CONCURRENCY - The end user is unaware of other users of the system. The current implementation does not support transactions. E.g. every change an end user performs on the system is immediately available to other users; there is currently no BEGIN/COMMIT/ROLLBACK for transactions.

- PARALLELISM - The end user is unaware of parallel execution of his requests.

- **Flexibility** - The system tries to be as flexible as possible with respect to performance. The three main points, which will give flexibility, are:

  - SERVICE CONCEPT - The system is designed around services: Each node has some services running, and these are responsible for certain requests (e.g. a data store service is responsible for storing/retrieving data). Those services are decoupled from each other and the only interfaces to the outside world are data items, which can be sent to it. Each service has to subscribe for keywords it is listening to (e.g. the data store service subscribes to the keywords "get", "put" and "delete").
  This allows extensions with new services, without having an impact on the design or structure of other parts of the system. (A detailed description can be found in section 4.1).

  - PLUGIN CONCEPT - Some services allow for further extension. E.g. the data store service allows for extensions with plugins to have different kinds of backend storage (file, database or others). The plugin structure has been designed to allow transparent access: The system itself does not know which exact plugin is in use; it just accesses it through a well-defined interface. (Detailed description can be found in section 4.2).

  - MESSAGE CONCEPT - The messages passed between the services are extensible as well. It defines a key-value concept, so each service can add new keys if it needs to. (A detailed description can be found in section 4.3).

- **Reliability / Availability / Redundancy** - For enhanced reliability and availability, the data can be replicated to other nodes. For write access (put, delete) the request is always routed to the master of the responsible data, and this master then transfers the data to its replicas. For read access (get), the requests are

load-shared among all master/replicas to enhance performance. To circumvent the problem of possible inconsistency, the routing is done at several levels, and if for some reason a request is routed to the wrong node, this node will relay the request to the correct one. In current implementation, no replication has been implemented yet, but the design allows for adding this functionality later.

- **Performance** - To gain optimal performance, all services are decoupled and run inside their own thread of execution. Together with other services, the system can run multi-threaded and can achieve performance gains on multi-processor machines. Also, the system is designed to allow asynchronous execution. This means all requests sent from one service to another one (local or remote) do not wait for an immediate response and continue to work with other data until the response eventually comes in. This is a very important aspect to allow high performance processing.

- **Scalability** - The scalability is highly dependent on the actual implementation of the communication/directory component (see section 4.4.3). From the design point of view, the system can scale to a very large number of nodes, as the system has been designed to hold as little status information as possible on the node itself (e.g. the algorithms do not store status information for other nodes).

- **Concurrency / Parallelism** - As the system is already distributed to multiple nodes, it already does parallel processing of data. To allow concurrent access to the system, there are two possibilities:

  1. The end user attaches his own application to a node that also stores data.
  2. The end user attaches his own application to a new node that does not store any data.

  In the former case, the node is a member of the distributed data structure and also a member of the end user application. In the latter case, the node is only a member of the end user application and connects to the system over the network. In both cases, the end user has concurrent access to its data.

- **Security** - No security counter measurements are currently in place. This implies that a malicious user could harm the whole system, e.g. by sending administrative messages, which could result in bringing down the whole system. However, the current design allows for extending the system with security measurements (such as encryption and authentication of messages as an extension to the sending/receiving entity).

Figure 24 shows the overall design with all its components (Services, Plugins, Messages with their queues and supporting components).

Figure 24: ERT System Overview

## 4.1   Service Concept

The system consists of several services, and each of them solves a well-defined part of the whole process. This allows for grouping functionality together and decoupling the system. It also allows for extending the system in the future with new functionality, which is not yet known or is application specific. Such extensions could be:

- A database application could implement a SQL [16] parser service, which would then use the underlying b-tree as a storage layer.

- A conversion service, which allows XML-based requests to be parsed and processed.

- A security service, which would grant/deny access on special criteria.

Generally such services only need be implemented if the standard user access, which only gives get/put/delete access to data, is not enough and access to lower-level parts are necessary to extend the functionality beyond the simple get/put/delete interface.

While services are not allowed to exchange data directly, there is a dispatcher, which will act as an intermediator between all services (see section 4.1.1 for details). The dispatcher decides where a specific request has to be forwarded; according to the supplied keywords, every service has to subscribe beforehand at the dispatcher. Using the service concept has the following advantages:

- **Decoupling** - Each service has a well-defined role and interface. A service can never call another service method directly. The interfaces between services are thread safe queues, and only the main component - the dispatcher - can put/get a workload from/to each service. Therefore, the only interface of a service is his queue, where he receives the workload from the dispatcher. All completed work is then sent again to the dispatcher, which then decides what to do with the work (e.g. send it to another service to process further).

- **Enhanced Performance** - Each service runs independently from the others. In the current implementation, each service runs within its own thread of execution, which scales with the number of CPUs of the underlying machine. Generally, it is not a good idea to have a large number of threads running at the same time, as each context switch of a thread will cause an overhead on the machine and will slow down the performance. However, with the current number of services, not many threads are running. Because the generation of the threads is done within the dispatcher, a further enhancement could be the introduction of a thread pool, where a predefined amount of threads are passed on to the busiest services.

- **Avoiding Threading Problems** - There is no limit how long a service needs to complete a certain task (e.g. the data store service takes longer to process, as it has to store data to disk, while other services have shorter processing time), so a multi-threaded system is necessary. Designing a well-functioning multi-threaded system is not an easy task (see section 2.5), as concurrent access to data is difficult to handle and the needed synchronization overhead will slow down performance.
  The service concept will solve this problem, as each service only runs within one single thread of execution. This way, the design and implementation of a service becomes easier, as it only has to be designed with a single threaded approach in mind, and no attention has to be paid to synchronization or concurrency issues. The only part of the system that is really multi-threaded is the dispatcher itself, as it will receive requests from all services and therefore has several threads of execution. But the dispatcher is a very small part of the system and consists only of a handful of lines of code and has been tested extensively (with respect to performance and accuracy).

- **Extensibility** - Each service has keywords assigned, which will act as routing indication for the dispatcher (e.g. the data store service will respond to "get", "put" and "delete" requests). This allows for easy extension, as new services just need to assign new keywords for their functionality and can be put in place without changing other services or the core system. Services could even be added/removed during runtime, just by registering/unregistering the keywords at the dispatcher.

### 4.1.1   Structure

Because every service has some basic thing in common, they all extend from an abstract service base (ServiceBase class). This class already implements features needed for message handling and registering service keywords.

It makes no sense to have several instances of the same service running, therefore every service needs to be implemented as a singleton [20, 26].

For a structural overview, see Figure 25.



Figure 25: Class Diagram for the Service Concept

An actual implementation of a service already gets the functionality of the service base class:

- **run** - Must be implemented in the actual service and must contain the run loop of the service. Typically, it consists of a *getWork* command and will then process the request. After completion, it will process the next request and so on.

- **registerAtDispatcher** - Must be implemented in the actual service. During initialization, this method will be called to register all its keywords at the dis-

patcher. It normally contains a list of *registerService* requests for all assigned keywords.

- **registerService** - With the *registerAtDispatcher* method, a list of those *registerService* requests will subscribe the service keywords at the dispatcher. It is not required to register the keywords during service initialization - this can be done at any time - but it is recommended to do it at the beginning, because if the dispatcher receives a request for an unknown keyword, it will just discard it.

- **getWork** - Will take the oldest request from the service queue. If the queue is empty, the thread will go into a waiting state until some work is pushed into that queue.

- **addWork** - Will put some work into the dispatcher queue. This method can be used if the work is of the "shoot-and-forget" kind and does not need to be handled as a transaction; therefore, no answer needs to be sent back.

- **transactionBegin** - Will put some work into the dispatcher queue. This method can be used if the work is a request that needs a response. If you use this method, the system makes sure that the response is routed back to this service. This also stays true if the request needs to be sent to another node. The response will always find its way back to its originating service.

- **transactionContinue** - Will put some work into the dispatcher queue. This method can be used if the request is received by this service but needs further processing by another service. You can do processing on this chunk of data and then pass it to another service. The response to this request will then be sent directly to the service, which originated the actual transaction.

- **transactionSplit** - Will put some work into the dispatcher queue. This method is a special case of a *transactionContinue* method. It will also relay the work, but is able to relay it to several services at once. This is useful if you need to spread some information to several nodes (e.g. if a "put" request needs to be stored on the master and also on all replicas, then the "put" request can be split and directly sent to all required nodes). A transaction monitor at the dispatcher of the initiating transaction will make sure that the final response of that transaction is collected from all sources and, if all sources report "successful operation", it will send back a "successful operation" message to the initiating service. If only one source sends back an "unsuccessful operation" message, the whole transaction is sent back as "unsuccessful" to the originating service.

- **transactionEnd** - Will put some work into the dispatcher queue. This method needs to be used if it is the final response to a request and the response should be sent back to the originating service.

- **getInstance** - This must be implemented in the actual service and should return the only instance of this service. The service has to be implemented as a

singelton [20, 26], and the *ServiceManager,* which will load the service, will use this method to get its service instance using standard JAVA introspection mechanisms. There is currently no other way in JAVA to force implementing a static method (such as forcing to implement a normal method using JAVA interfaces). In any case, you would get a runtime exception if the *ServiceManager* tried to load this service but could not find this *getInstance* method.

The complete processing of requests from one service to another is shown in the following Figure 26



Figure 26: Processing of a Request

1. Service A sends a request to the dispatcher (either as *addWork*, *transactionBegin*, *transactionContinu*e or *transactionEnd*).

2. Dispatcher will pick up the oldest entry of its queue and start processing it.

3. Dispatcher will lookup in its registered keywords list to find where it has to put the request.

4. Dispatcher will find a suitable Service B and puts the request in its queue.

5. Service B will pick up the oldest entry of its queue and start processing it.

### 4.1.2   Usage

- **Creation** - The instance is created by the *ServiceManager*, which manages all services that should be loaded during application startup. This *ServiceManager* will use the *getInstance* method to get the single instance of this service.

```
static public synchronized DataStoreService getInstance() {
  if (null == instance) {
    instance = new DataStoreService();
    instance.dataStoreManager = DataStoreManager.getInstance();
    logger.info("DataStoreService started");
  }
  return instance;
}
```

Figure 27: Creation of a Service

- **Registering** - The *ServiceManager* uses *registerAtDispatcher* to register all keywords for this service at the dispatcher. The service is then set up to receive requests for those keywords. There is one special keyword, "default", which can be implemented by one and only one service, which will then be chosen by the dispatcher, if no service for a particular keyword has been found.

```
public void registerAtDispatcher() {
  registerService("putdb");
  registerService("getdb");
  registerService("getnextdb");
  registerService("getprevdb");
  registerService("deletedb");
  registerService("getsizedb");
}
```

Figure 28: Registering at Dispatcher

- **Run Loop** - After the initialization phase, the *ServiceManager* will execute the run loop of the service, and the service will do the actual work.

```
public void run() {
  Thread.currentThread().setName("ERT.KeyPartitionService");
  while (true) {
    DataItem dataItem = getWork();
    String operation = dataItem.getOperationType();
    if (operation.equals("get")) {
      keyPartitionManager.get(dataItem);
    }
    ...
  }
}
```

Figure 29: Run Loop

- **Responding** - After or during the processing, one or more new requests/responses are sent to the dispatcher. This can be done inside the service itself, or within subcomponents of the actual service.

```
DataItem response = dataStoreManager.get(dataItem);
transactionEnd(response);
```

Figure 30: Responding

## 4.2   Plugin Concept

Besides extending the functionality with new services, some core services have the possibility to change their functionality. Currently there are two services which have a common plugin structure: *NodeCommunicationService* and *DataStoreService.* Extensions to services could be:

- New routing mechanisms for messages (e.g. Zeroconf based [35], P2P-overlay based [15, 17])

- New messaging methods for communication between system instances (e.g. JAVA Remote-Method-Invocation, XML based, ASN.1 based [6, 45])

- New storage mechanisms for actual data (e.g. file based, database based)

The main idea behind the plugin concept is to allow for exchanging parts in an already implemented service with other plugin components. Instead of building a complete new service from scratch, extending existing services with plugins reduces the development time. Using the plugin concept has the following advantages:

- **Reducing Development Efforts** - By re-using the already available service component with all its business logic, the time to develop and test a new service can be reduced, as only part of the plugin itself needs to be written and

tested. The plugin itself needs to implement a JAVA interface and has to con-
form to its method declaration. The plugin will then be loaded by the service
itself as a delegate/proxy (see delegate/proxy pattern in [19, 26]). This means
the service or the whole system does not know which particular plugin is in
use.
Also, the plugin can use the already proven base classes, such as logging, con-
currency packages and data structures.

- **Extensibility** - As you would use the service concept if you wanted to make
  structural changes to the whole system (e.g.  by introducing new keywords,
  functionality), you can use the plugin concept to make smaller changes to al-
  ready developed services.

- **Abstraction** - As the system itself is not aware of the actual plugin implemen-
  tation, the whole system will benefit by the reduction of complexity that needs
  to be handled.

### 4.2.1   Structure

The actual structure is dependent on the service and the interface it provides for
plugins, so the number and signature of the methods can change for each service.
For a structural overview, see Figure 31.



Figure 31: Class Diagram for the Plugin Concept

- **method1 / method2** - The interface specifies which methods should be part
  of the *Plugin*.  Those methods need to be implemented at the *Plugin*.  The

*PluginDelegate* will also implement those methods, but only as an envelope to the actual plugin. This means the implementation of the *PluginDelegate* will use the same method in the *Plugin* instead and therefore acts as a proxy (see section 4.2.2 for details).

- **PluginDelegate** - During instantiation of the *PluginDelegate*, it will load the actual *Plugin* using JAVA introspection mechanisms. In the current design, the *Plugin* needs to have the same name as the *PluginDelegate*, and it needs to reside inside its own sub package (see JAVA code for examples).
  The *Service* will only access the *PluginDelegate* to actually reach a plugin. This abstracts the actual plugin implementation from the rest of the system.

### 4.2.2   Usage

- **Creation** - The *PluginDelegate* needs to load the actual implementation of the *Plugin*. This is done in the constructor of the delegate, as shown below. The delegate is then initialized (e.g. by setting properties, which are needed for setup).

```
static private String plugin = "berkeleyje";

public DataStore(String instanceId) {
  try {
    delegate = (DataStoreInterface)
                   (Class.forName("ch.nix.ert.store."
                   + plugin + ".DataStore").newInstance());
    delegate.setDatabaseId(instanceId);
  } catch (InstantiationException e) {
    e.printStackTrace();
  } catch (IllegalAccessException e) {
    e.printStackTrace();
  } catch (ClassNotFoundException e) {
        e.printStackTrace();
  }
}
```

Figure 32: Creation

- **Delegating** - If a method is accessed, the delegate acts as a proxy and forwards the request to the actual *Plugin*.

```
public DataItem getDataItem(DataItem data) {
        return delegate.getDataItem(data);
}
```

Figure 33: Delegating

- **Implementation** - The actual implementation of the functionality inside the

*Plugin.*

```
public DataItem getDataItem(DataItem data) {
  //extract key from data
  long key = data.getUserKey().getId();
  //check, if connection to database is already open
  if (myDatabase == null) {
    this.openDatabase();
  }
  try {
    // Create a pair of DatabaseEntry objects. theKey
    // is used to perform the search. theData is used
    // to store the data returned by the get() operation.
    DatabaseEntry theKey = new DatabaseEntry();
    DatabaseEntry theData = new DatabaseEntry();
...
```

Figure 34: Implementation

## 4.3  Message Concept

The current architecture is based on message passing. All services pass messages along to other services. Those can be on the same node or on a remote node. The basic idea of this message concept is to introduce a common way for intercommunication between services, which also allows for future extension. Such extensions could be:

- New yet unknown services need specific data, which needs to be transported over the already deployed network.

- In addition to the raw key-value data from users, other properties such as permission-related attributes and metadata might be introduced.

- New communication protocols (e.g. switching to XML, ASN.1 [6, 45]) could be introduced; this should not affect the implementation of the services.

As the future is unknown, the messaging should be decoupled as much as possible from the actual implemented services, so additional attributes could be added without affecting already developed services (backward, forward compatibility [75, 79]). Generally, using the message concept has the following advantages:

- **Backward Compatible** - New versions of ERT should understand old messages, as the messages themselves are based on a key-value mechanism. In newer versions, new keywords can be introduced, but the old keywords are still recognized. A system with a mixture of old and new systems is possible and in the case of a rolling update (an update with zero downtime) unavoidable.

- **Forward Compatible** - Old versions of ERT should understand new messages, again because the old key-values are still present in the new messages. Of

course, this has the limitation that features introduced in the new version are not available on the old system. However, because the message concept will still understand the grammar of the message, it will not alter the message itself, but just not react to new attributes. This is also important during a rolling update of a system, where it is unavoidable that at a certain time, multiple versions of the system are running at the same time.

- **Extensibility** - Because of the key-value nature (See section 4.3.1 for details), new attributes can be introduced at any time. Even when new attributes are introduced, the communication component is still capable of sending the messages, even if they contain new/unknown attributes. The same is true for services: Each service only looks at those attributes it is interested in and will not fail if new attributes are introduced.

- **Abstraction** - The structure of the actual message is hidden for the developer; he does not know how the message is built. The message itself is a lightweight object, which basically contains an array. The keys are integers and the values can be of any kind. Making the values as lightweight as possible, e.g. using only strings, integers or simple objects to speed up the processing speed, is encouraged. This is important, as each sending and storing activity will serialize the message object and serializing complex object structures can be slow.

### 4.3.1   Structure

The message concept is based on a class named *DataItem*, which consists of an array and the idea that every attribute is represented by a key-value pair:

- **Key** - The key is the representation of the attribute. It is of type *int* (integer) and in current implementation all keys are on the same level. There is no namespace. Every attribute is identified by its unique key. In the future, as the attributes might grow, a nested approach in which keys are grouped and nested together could be necessary.

- **Value** - The value is the attribute content. It is of type *Object* and therefore can be any kind of information. Having only simple objects is encouraged, as there is a performance penalty to convert complex objects to the actual binary representation needed for storing and communication between nodes.

To make it easier for the developer, all key-value attributes have accessor methods, which help setting/getting/deleting attributes. It is favorable to use these accessor methods, because using the array directly will always result in a return value of type *Object,* which would need to be cast to the needed object by hand afterwards. Using accessor methods does this for you.

For a structural overview, see Figure 35

Figure 35: Class Diagram for the Message Concept

- **Serializable** - There is no method to implement for the *Serializable* interface. Nevertheless, *DataItem* needs to implement it, so this class can be automatically serialized for use in data store and communication layer as needed.

- **clone** - This allows cloning the object. Sometimes this is needed if you want to split the transaction (see section 4.1.1 for details) into multiple pieces and every piece contains the same *DataItem* (e.g. to replicate the very same data to several nodes). It is also needed if multiple references of the same *DataItem* are sent around inside the same node, as then multiple threads might modify this *DataItem* at the same time. Cloning can help in these circumstances. Generally, the *DataItem* is implemented in a thread-safe way, so even multiple threads could eventually modify the same *DataItem*. Nevertheless, some services may not allow the *DataItem* change at the same time by another service and therefore need to implement cloning.

- **getValue** - Gets a value for a given key. The key is always of type *int* and this method returns an object of type *Object*. Using accessor methods for the keys instead is encouraged, as they provide optimized versions on accessing the actual value.

- **putValue** - Assigns a value for a given key. The key is always of type *int* and value of type *Object*. Using accessor methods for the keys instead is encouraged, as they provide optimized versions on storing the value.

- **removeValue** - Removes a given key. The key is always of type *int*. Using accessor methods for most of the keys instead is encouraged.

- **existValue** - Checks for the existence of a certain key. The key is always of type *int* and will return a *boolean* (false/true).

- **getClone** - This is a more convenient method to clone a *DataItem*, as it already returns an object of type *DataItem* (as clone itself only returns a generic *Object*).

- **get / set / remove** - These are accessor methods for the most common used keys. Currently all used keys have such accessor methods and using them is strongly encouraged, as they already provide basic conversion to the final object type.

### 4.3.2   Usage

- **Getting Values** - There are two ways of getting values: One is by using the general method, and the other is by using the accessor methods.

```
Node originator = (Node) dataItem.getValue(
          OperationInterface.DIRECTORIGINATOR);
```

Figure 36: Getting Values - General

```
Node originator = dataItem.getDirectOriginator();
```

Figure 37: Getting Values - Accessor Method

- **Setting Values** - There are two ways of setting values: One is by using the general method, and the other is by using the accessor methods.

```
Node node = new Node();
dataItem.putValue(
          OperationInterface.DIRECTORIGINATOR, node);
```

Figure 38: Setting Values - General

```
Node node = new Node();
dataItem.setDirectOriginator(node);
```

Figure 39: Setting Values - Accessor Method

- **Removing Values** - There are two ways of removing/deleting values: One is by using the general method, and the other is by using the accessor methods.

```
dataItem.removeValue(OperationInterface.DIRECTORIGINATOR);
```

Figure 40: Removing Values - General

```
dataItem.removeDirectOriginator();
```

Figure 41: Removing Values - Accessor Method

## 4.4   Components

This section will describe the most important components, which form the ERT system.

### 4.4.1   Key Partition Service

This service is responsible for deciding if this node can handle an incoming request or if it should be relayed to another node. Furthermore, it handles all data migration requests and in future versions of ERT will also handle replication. In current implementation, the *Key Partition Service* listens to the following commands:

- **put** - Will first check to determine whether this node is responsible for this key-value pair. If so, it will store the key-value by using the *Data Store Service* (see section 4.4.2) method: putdb. If the *Key Partition Service* decides it is not responsible for this key-value pair, it will try to relay the request to the node it thinks is responsible for it.

- **get** - Will first check to determine whether this node is responsible for this request. If so, it will retrieve the value for this key by using the *Data Store Service* (see section 4.4.2) method: getdb. If the *Key Partition Service* decides it is not responsible for this key-value pair, it will try to relay the request to the node it thinks is responsible for it.

- **getnext** - Will first check to determine whether this node is responsible for this request. If so, it will retrieve the value for the next key by using the *Data Store Service* (see section 4.4.2) method: getnext. If the *Key Partition Service* decides it is not responsible for this key-value pair, it will try to relay the request to the node it thinks is responsible for it.

- **getprev** - Will first check to determine whether this node is responsible for this request. If so, it will retrieve the value for the previous key by using the

*Data Store Service* (see section 4.4.2) method: getprev. If the *Key Partition Service* decides it is not responsible for this key-value pair, it will try to relay the request to the node it thinks is responsible for it.

- **delete** - Will first check to determine whether this node is responsible for this request. If so, it will delete the corresponding key-value pair by using the *Data Store Service* (see section 4.4.2) method: deletedb. If the *Key Partition Service* decides it is not responsible for this key-value pair, it will try to relay the request to the node it thinks is responsible for it.

- **nodeload** - Will return the current load of this node. The return value is a number between 0 and 1. 0 means: Node is idle and 1 means: Node is under heavy load.

- **firstnode** - During startup phase, this command indicates that this node is the first node on the whole network. This implies that this node does not try to probe for other nodes during startup and will assign the whole data range (all possible keys) to itself.

- **join** - During startup phase, this command indicates that this node should be part of an already existing network. It will then probe some random nodes to find out which node currently has the highest load, and will then ask this node to migrate some data from it.

- **migrationprobe** - Is similar to the join command. During runtime, the system will probe its neighboring nodes (those nodes which are responsible for the next higher or lower data range) every few minutes to see if they have comparable load. If it finds that the neighbors are more heavy loaded compared to its own load factor, it will start a migration. This way it can be assured that the overall load is evenly distributed among all participating system nodes.

- **migrationstart** - This request will start a migration (see section 4.5.4). As each node can only be part of one migration at a time, it will first check to determine whether this node is already doing another migration. If not, it will respond with the data range to migrate.

- **getmigration** - This request will actually transfer the data from one node to the other.

- **migrationend** - If the migration is finished and all necessary data has been transferred to the other node, this request will indicate that the data has been successfully migrated. The two participating nodes will then adjust their *Node Directories* and carry out other administrative actions.

**4.4.1.1   Responsibility**   The *Key Partition Service* will always check to determine whether the local node can handle the incoming request. If not, it tries to relay the request to another node, which it thinks is responsible for. Of course, the newly

chosen node again checks the responsibility and will further relay the request until the correct node is found. Under normal circumstances, in which the node directory is up-to-date, no relaying is necessary. Relaying only happens if one node directory has temporary outdated information.

**4.4.1.2   Migration**   The *Key Partition Service* is responsible for data balancing and will migrate the data automatically if necessary. During startup of a new node, it will probe a few random nodes to find out which is the most loaded node and will then start migrating some data to the new node. After this first migration (the new node joins the network), the nodes will continuously watch their neighboring nodes, and, if their load starts to climb, request a migration to relieve overloaded nodes. See section 4.5.4 for more details.

**4.4.1.3   Replication**   In current version of ERT, no replication feature has been implemented. Replication already has been prepared in the design of the application and should be straightforward to implement, as the logic should be about the same as for migration.

### 4.4.2   Data Store Service

This service is responsible for actually storing the data permanently, e.g. to disk. It uses the *Plugin Concept* described above to allow replacing the actual data store with another implementation. The following commands are supported:

- **putdb** - Will store a tuple (key-value) to the storage. The outcome of this operation will be sent back by returning true/false.

- **putdbnoack** - Will store a tuple (key-value) to the storage. There is no response sent back, so there is no guarantee that the tuple was successfully stored.

- **getdb** - Will retrieve a tuple (key-value) from the storage. The response will contain the complete key-value pair.

- **getnextdb** - For a given key, it will respond with the next available key-value pair. E.g. if your database contains the following key-value pairs: (1→DATA1), (4→DATA4), (12→DATA12), a *getnextdb* request for key 1 will return the tuple (4→DATA4).

- **getprevdb** - For a given key, it will respond with the previous available key-value pair. E.g. if your storage contains the following key-value pairs: (1→DATA1), (4→DATA4), (12→DATA12), a *getprevdb* request for key 4 will return the tuple (1→DATA1).

- **deletedb** - Will delete a tuple (key-value) from the storage. The outcome of this operation will be sent back by returning true/false.

- **getsizedb** - Will return the current amount of tuples (key-value pairs) of the datastore.

In the current implementation, only a BERKELEYDB JAVA EDITION [65] storage is available, but other plugins can be built and used instead.

**4.4.2.1   BERKELEYDB**   BERKELEYDB [65] is an implementation of a balanced tree (b-tree) and is available from Sleepycat Software. Two different implementations are available: BERKELEYDB and BERKELEYDB JAVA EDITION. The former is written in C-language and is a proven product, used in multiple commercial and non-commercial applications. There are JAVA bindings, through the JAVA Native Interface (JNI), so JAVA applications can use the C version of this implementation. The newer BERKELEYDB JAVA EDITION is a re-write of the whole data-store in JAVA language. There is no need to install any C libraries anymore in the system, as this version is now 100% JAVA.

Some general features of BERKELEYDB are (taken from [65]):

- **Open Source** - The sources for BERKELEYDB are available for download and BERKELEYDB can be integrated into your own open source software for free, although if you want to include it in your own commercial applications (where you don't want to make the sources available), a license has to be bought.

- **Small Footprint** - In contrast to full-blown databases (e.g. Oracle, DB2), BERKELEYDB only needs very limited space on disk for the executable, and the memory footprint of the loaded executable is very low.

- **Thread-Safe Library** - BERKELEYDB is completely thread-safe; therefore, multiple threads can be read/written to the database at the same time. This feature is currently not used, as this *Data Store Service* only runs within a single thread.

- **File System Integration** - The database itself is stored on a normal file system. There is no need for a special partitioned hard disk or other storage medium.

- **Large Databases** - BERKELEYDB supports databases up to a size of 256 terabytes and key-values up to 4 gigabytes in size.

- **Transaction** - BERKELEYDB allows transaction. This means that several operations can be performed on the database and they will only be visible to others if the whole transaction is executed. Otherwise, BERKELEYDB will roll back to the state it was in before the transaction started.

- **Check Points** - If BERKELEYDB or its hosting application crashes, BERKELEYDB automatically brings the database into a consistent state after the next startup.

- **Hot Backup** - A backup of the database can be done during normal operation. There is no need to shut down the application to do a backup of the data.

### 4.4.3   Node Communication / Directory Service

This service is responsible for sending requests to other nodes and consists of two main components:

1. **Node Communication** - This part is responsible for sending direct requests to nodes. Such requests can be administrative in nature (e.g. status messages) or application related (e.g. get/put/delete requests). The general dispatcher (see section 4.1) tries to route the messages, and if he can't find the required service on the local node, he will put it on the *Node Communication Service* and this service will then send out the request to the correct node.

2. **Node Directory** - As the *Node Communication Service* needs to know which node is responsible for which data range, the *Node Directory* maintains a list (this list can be real or virtual; e.g. for smaller systems, the list could be stored on each node. For larger systems, the list can be distributed over all or parts of the nodes) and just responds to requests like "give me the IP of the node which is responsible for key a".

Both components are implemented following the *Plugin Concept*. Depending on the needs, both components (communication and directory) could be implemented as one single large component if needed, or separated.

The following commands are supported:

- **send** - Will send out this data item. Depending on the routing indication of this data item, the service will first do a lookup in its *Node Directory* or if an IP is already supplied in the data item, it will send it directly to the corresponding node.

- **default** - The dispatcher (see section 4.1) will automatically send all requests for which it does not find a local responsible, to a service which subscribes the "default" keyword. The *Node Communication Service* will then receive all requests that cannot be handled by the local instance and therefore need to be transferred to another node.

- **publish** - If there is a new key-group to be added to the local instance (so the local instance will get a new data range assigned), the *Node Directory* needs to be called with this command, to update its directory. It will then pass this information to its *Node Directory* plugin, which will carry out the necessary steps to publish this information, e.g. to send a broadcast to all other nodes, telling them that this node is now responsible for this new data range.

- **unpublish** - If the current node is no longer responsible for a certain data range, it needs to send this message to the *Node Directory*. It will then pass on this information to its *Node Directory* plugin, which will carry out the necessary steps to unpublish this information.

In the current implementation, two communication plugins (RMI, NETTY2) and one directory plugin (JMDNS) are available.

**4.4.3.1    RMI**    This communication implementation on top of the standard JAVA
Remote Method Invocation (RMI) [67] is a rather slow but straightforward method
for sending messages from one node to another. It basically uses two main parts for
sending messages:

- **Serialization** - A RMI message is basically a full JAVA object that is being sent
  over the network. Before the sending can be done, RMI does a serialization
  of the object into a byte stream. That means all relevant data of an object (in-
  stance variables, static variables) are converted to a flat byte stream, including
  administrative data to de-serialize the data stream afterwards. RMI automat-
  ically uses an internal serialization mechanism, so the application developer
  does not need to write its own. He can, however, by implementing special
  serialization methods (e.g. for better serialization performance).



Figure 42: Serialization, De-Serialization

- **Sending** - After a byte stream has been generated through serialization, RMI
  tries to deliver it over the network.

The main use of this plugin is to have a running communication layer very quickly,
as most implementation details are already included in the RMI library. As a down-
side, the performance is quite low. A test between two machines showed that this
mechanism allows around 80 messages per second to be sent over the network. The
main bottleneck is the serialization mechanism, which takes the most time.

**4.4.3.2    NETTY2**    This communication implementation on top of NETTY2 [46]
uses the new JAVA NIO-package (New I/O) [68], which claims to support superior
network performance. The main difference between the old networking libraries and
the new are: The old was using one thread for each listening network socket, while
the new is only using a user-specified number of threads for all network sockets.
This is especially useful for servers, which need to listen to many concurrent network
channels from different nodes.

Figure 43: Old vs. New JAVA Networking

The main advantage of this plugin is its better performance. A test between two machines showed that this mechanism allows around 700 messages per second to be sent over the network (around 10 times the performance of RMI). As NETTY2 or NIO are only libraries for doing raw network transfer, the downside is that the application developer needs to build its own serialization/de-serialization mechanism. This has been done for ERT, which implements an efficient way to encode/decode messages in a way, similar to ASN.1 [6, 45].

**4.4.3.3    JMDNS**    This is a directory implementation on top of JMDNS [71], which conforms to the Zeroconf/Rendezvous [2, 35] specification.



Figure 44: JMDNS Message Flow

It is responsible for maintaining the directory of ERT and holds a list, which part of the data is stored on which node. It implements a broadcasting mechanism to transmit changes to its directory to other nodes. As changes to the data range do not occur very often, this mechanism does not need to be very efficient. Also, the

problem of propagating time is solved, as the system works even if one node has outdated information about the other nodes. In this case, the lookup just takes a bit longer, as the node that received the request and does not feel responsible will just relay it to the correct one.

### 4.4.4   Callback Service

This service is responsible for automatic callbacks within the ERT application. Normally, if you need periodic execution of a code (e.g. to collect statistic information every minute, clean-up), a separate thread that will execute the code from time to time is introduced. This can lead to poor performance, as every such situation would need a separate thread (e.g. if you need to update statistical information every minute, clean up every 10 minutes and poll a server every 5 seconds). *Callback Service* tries to outsource this mechanism to its own service. Additionally the timely execution is always done over the normal service interface, so the callback service never executes code directly; it just puts the callback into the service queue, which requested the callback before. This way, every service can still run within one single thread (see section 4.1), and no synchronization overhead needs to be put in place.

There are two types of callbacks:

- **calllater** - The callback is executed only once after a specified amount of time.

- **callcontinous** - The callback is executed at every specified time (e.g. every minute).



Figure 45: Callback Service

1. Service requests a callback (e.g. a continuous callback)

2. After the time is up, the *Callback Service* puts the original request into the Service Queue[9]

---

[9]For simplicity, the role of the dispatcher is not shown here. All requests for *Callback Service* or *Service A* normally go through the dispatcher.

### 4.4.5 User Service

This service is part of the user application; the user of ERT will interact only with this *User Service*. It supports wrapper methods for commonly used functions and will hide the complexity of a distributed system. The methods provided by the *User Service* are:

- **put(key, value)** - Will store the key-value pair and will not block further execution. Whether the request can be executed successfully or not, it will be returned at a later stage using an asynchronous response.

- **putSync(key, value)** - Will store the key-value pair and will block until a response (successful/unsuccessful) is returned by the ERT system.

- **get(key)** - Will retrieve a value for this given key. The request is executed asynchronously, and the response is sent back at a later stage.

- **getSync(key)** - Will retrieve a value for this given key and will block until a response (the value or a status message indicating a failure) is returned.

- **getNext(key)** - Will retrieve the next value for this given key. The request is executed asynchronously, and the response is sent back at a later stage.

- **getNextSync(key)** - Will retrieve the next value for this given key and will block until a response (the value or a status message indicating a failure) is returned.

- **getPrev(key)** - Will retrieve the previous value for this given key. The request is executed asynchronously, and the response is sent back at a later stage.

- **getPrevSync(key)** - Will retrieve the previous value for this given key and will block until a response (the value or a status message indicating a failure) is returned.

- **delete(key)** -Will delete the value associated with this key and will not block further execution. Whether the request can be executed successfully or not, it will be returned at a later stage using an asynchronous response.

- **deleteSync(key)** - Will delete the value associated with this key and will block until a response (successful/unsuccessful) is returned by the ERT system.

The only interface between the user application and ERT will be the interface described above, so a user application never uses other services or other parts of the ERT system.

### 4.4.6 Main Service

The only responsibility of the *Main Service* is to initialize the system and listen to application requests (such as shutdown, startup, dumping version information).

## 4.5    Message Flow

This section will describe the most common message flows of ERT. For simplicity, the role of the dispatcher is omitted.

### 4.5.1    Get Requests



Figure 46: Get Request

The *UserService* will initiate a get request (*TransactionBegin*) and if the request can't already be answered at the current node (e.g. because the node does not contain any data, as it is not part of the storing network and only acts as a user service, or the data is not part of its own data partition), the *NodeCommunicationService* will find the correct node which is the holder of the requested data and will send the request there. If for some reason the *NodeCommunicationService* has outdated information and therefore sends the request to the wrong node, the receiving node will relay the message to the correct one.

The *DataStoreService* will answer the request and end the transaction (*TransactionEnd*) and the transaction mechanism will then route the response to the node that originated it.

Additionally there are other similar get requests not shown in the message flow chart. There are also the *getNext/getPrev* commands, which work similarly to the normal *get* command.

### 4.5.2   Put Requests

Figure 47: Put Request

Put requests use the same procedure as get requests. A put request normally has no return value. In this architecture, a put request will always return, whether the request has been successfully executed or not. This is important, because in a highly distributed system with unpredictable network latency, it is not certain that the put request can be processed and stored successfully.

### 4.5.3   Delete Requests

Figure 48: Delete Request

Works similarly to the put request above. It will also return whether the delete request can be executed successfully or not.

### 4.5.4    Migration Requests



Figure 49: Migration Request

Migration has the most complex message flow and was designed to be as fault tolerant as possible. In addition to the mechanism pictured above, different timeout mechanisms have been put in place so the whole migration can be rolled back if something happens. It is interesting to note that in case of a migration roll-back, no data have to be moved again, as during migration no data is being deleted at the original node and all write/delete requests during migration update both parties simultaneously.

The whole migration process starts by probing the neighboring node at certain time intervals with a node load request. If the node load of the neighboring node is much higher then the local load factor, the node that has initiated the node load request will start a migration request. The other party will respond with a migration request response, which contains the exact amount of data that could be migrated, including the lower/higher boundary of the data to migrate (e.g. range $[p, t]$). The migration initiator will then check the range and choose a sub/whole-range of it and start requesting the data (*getMigration*). Once all data have been transferred, the initiator will end the transaction by sending a *migrationEnd* message.

Note that there is no deletion of data. Even if the migration has been successfully terminated, the data are not deleted immediately. There is a background task running on each node, which will take care of deleting all unnecessary data after some time.

### 4.5.5   Administrative Requests



Figure 50: Administrative Requests

The only administrative requests at the moment are requests to join a running network, and for the first node, there is a request to tell the node that it is the first node on the whole network.

# 5   Evaluation

This section will provide an overview of the tests done with the prototype and will list some issues found or things which have worked out well.

## 5.1   Benchmark

Benchmarking has been done at the student computer lab at the IFI (Institute for Informatics, University of Zurich), running on multiple machines. All machines have the same hardware and software setup (see CD-ROM for a detailed machine profile, including all installed software and hardware).

- **Hardware** - PowerMac G4, dual CPU, each running at 450 MHz with 256-640 MB of RAM.

- **Software** - MacOSX 10.3.5 running JAVA version 1.4.2_05-141.3 and HotSpot client VM 1.4.2-38. ERT was running in version 165 and was configured to use the NETTY2 (see section 4.4.3.2) communication layer.

- **Network** - All machines are connected using switches on a 100 Mbit/s (Fast Ethernet) network.

To perform the tests, the machines are divided into three roles:

- **Collector** - A single machine, who will collect all the data received from the load generators and will calculate the overall load.

- **Load Generator** - Those machines generate the load to the test machines. They run the same ERT software, but do not perform any storage on their own. They just send out requests, e.g. a *put(key, value)*, and wait for their responses. The number of responses per second is sent to the collector every 10 seconds. This figure is used to calculate the overall performance of the whole system. It is important to note that not the amount of requests sent out per second is recorded (as this is much higher), but only the amount of responses returned.

- **Test Machine** - These machines process the requests received from the load generator and send back the result. There is no measuring of the processed number of requests, as this is done solely at the load generators.



Figure 51: Benchmark Setup: Collector (C), Load Generators (LG) and Test Machines (T)

In the benchmark printed below, 1 collector, 20 load generators and up to 42 test machines are used. The load of all 20 load generators is added up to get the overall load.

Execution of the benchmark has been automated and is always done in the same manner:

1. The database on each node is cleared.

2. ERT is started on all test machines.

3. Waiting until the test machines are stabilized and all data-ranges have been distributed equally to all nodes.

4. Collector is started and is ready to receive results.

5. Load generators are all started simultaneously and start sending requests. To prevent overloading the system, the load generators start with a small load and will increase their load every second, until the maximum possible throughput is reached.

6. This setup runs for a few minutes and is then stopped.

7. The results from the collector are noted down. Always the highest three num-
   bers in a row are averaged and used for plotting the data below.



Figure 52: Benchmark: Random put, 10 bytes of data



Figure 53: Benchmark: Random put, 1024 bytes of data

As shown in the benchmark result, ERT will scale pretty well as the number of nodes increases - although it did not scale entirely linear, this could be related to network stalling, as all participating machines share the same physical network connection.

It should be noted, that each transaction consists of two messages being sent: One request and one response, so the number of messages processed is twice the number as shown above.

Benchmarking was only done for "put" requests and not for the others like "get", "delete", as all requests have the same characteristics (all transactions are based on two messages being sent, see section 4.5), with the only difference that the destination node will either store, retrieve or delete data from its internal database. The "put" request is the most time consuming request anyway, as the underlying storage architecture (BERKELEYDB [65]) needs to allocate new disk-space and handles all insertions into its internal b-tree. Beside that, in current testing setup, all "put" requests have been executed in a random manner, to retrieve those values afterwards a lot of "get" requests would return an empty result, if no key-value pair for this particular "get" request is found.

## 5.2   Strengths / Weaknesses

During testing and developing, the following strengths and weaknesses were identified:

- **Strengths**

    - The complexity of adding new nodes to the network is very low, as there is not much configuration to be done. Basically, there are two types of nodes: The first one creates a new ERT network and needs to be started with a "first" flag. Then all other nodes can be added by supplying a "join" flag, and will then automatically try to join the existing network and assume responsibility for their own data range.

    - The overall performance of the system is good, although no intensive performance tuning has been done so far.

    - The system uses both CPUs efficiently, as the load of the system climbed almost to the maximum of 200% CPU utilization.

- **Weaknesses**

    - A client who wants to access the distributed data needs to open a network connection to all those nodes where the specific data is located. If a client needs to access all available data (e.g. by randomly accessing all of them), it potentially will open a network connection to all participating nodes on the network. This can cause problems if the network is very large (e.g. thousands of nodes).

– The new JAVA networking library (NIO), which is used in the NETTY2
plugin, as described in section 4.4.3.2, offers great performance, but it
was added to the JAVA library only very recently, as of JAVA version
1.4. It seems not all vendors of JAVA have a robust implementation of
it; e.g. the JAVA version included in Apple's MacOSX suffers from sta-
bility problems, which can cause the whole machine to crash. There is
a workaround, but this decreases overall performance a bit. A bug re-
port has been filed at Apple (Bug-ID 3866318), so this will hopefully be
resolved later.

– Transferring large amounts of data, which can happen if the key-value
pairs are very large in size, will decrease the overall performance of the
system. As in a distributed system, the networking performance is usually
the bottleneck. Sending large chunks of data will occupy the network, so
other concurrent access to it will suffer. Also, sending large pieces of
data will increase the memory footprint of the application, as all the data
needs to be kept in the memory at some time.

– Migration of data under load is currently not very efficient and still has
some bugs, which was an additional shortcoming during benchmarking.

# 6 Conclusion

This section will provide a short summary of the work and give an outlook.

## 6.1 Introduction

Current problems in data storage are the growing amount of data, the increasing level of access to those data and the need for finding the correct data. This leads to the fact that today's systems are built on distributed foundations and a trend towards more structured data. While there are some solutions for distributing data, most of them are done at a very high level.

This work will introduce a new low-level approach to distributing data - ERT (Extended Root Tree) - which is based on a distributed sorted key-value architecture. It can be used as a building block for higher-level applications as a storage tier.

## 6.2 Distributed Data Storage

In today's world, we see tendencies that indicate the amount of data is growing. More than 5 terabytes of data are produced every year. Moreover, companies, as well as private individuals, have data distributed to many places (e.g. company computer and personal laptop); so finding data is becoming increasingly difficult. A solution for finding the correct data is to structure it using several mechanisms. Depending on the use-case and application, low-level structures, such as key-value based algorithms, are used, or more high-level approaches, such as storing the data into databases, are chosen.

To discuss distributed data storage, we first need to define what characteristics a general distributed system has and then build a system, which should fulfill those requirements.

For a better understanding of how others have built distributed data storage systems, an overview of the large database vendors (Oracle, IBM, Microsoft) is given, and it shows that they all do the data distribution at a very high level. Further general research projects that focus on low-level distribution of data are discussed, and there we find two main directions: Smaller systems allow for sorted key-value access but most larger systems can only handle unsorted key-value access.

Regardless, for monolithic or distributed systems, they need to store data. As each machine only has direct access to its own hardware, the actual low-level storing is always done locally and a brief overview of such low-level storage is given. It shows that most systems implement a balanced tree (b-tree)-like structure for low-level storage.

Another problem in distributed systems is performance aspects: As distributed systems need communication between them, optimizing the network traffic is a key element of a successful system. Another one is asynchronous processing of requests to enable higher overall throughput of the system.

## 6.3  Overview of ERT

The main idea behind ERT is to split the data range into sub ranges and assign each computer on the network one or multiple of such sub ranges. The administration of the sub ranges (e.g. where to find a specific sub range) is done inside a distributed directory. The ERT approach does not favor a specific directory, so depending on the needs and size of the system, different such directory implementations can be used. For smaller systems, static directories could be used, for medium-sized systems, broadcast mechanisms and for very large systems, P2P (Peer-To-Peer)-like systems. The actual storing of key-value pairs on each node is done using a conventional balanced tree. Also, here ERT does not favor a specific implementation. Again, depending on the needs, different such balanced tree-like structures can be used.

## 6.4  Implementation of ERT

The prototype is being implemented in JAVA and basically consists of three main concepts: The service concept, which is responsible for decoupling the whole system into subcomponents, the plugin concept, which is responsible for abstracting the actual implementation from the services and the messaging concept, which is responsible for all communication between and inside all components.

## 6.5  Evaluation

Benchmarking showed that the system scales well as the number of nodes grow. There are bottlenecks, however, as networking bandwidth can limit scalability.

## 6.6  Outlook

The current prototype is suitable for distributed low-level data storage and is fast enough depending on the use case. Testing on how big the system can scale is still needed. The network interface could then become a bottleneck limiting further scalability. Compared to a monolithic system, the performance improvements in this distributed data structures only has an effect after an initial size is reached. Therefore, a system consisting of only two machines is slower than a monolithic system.

While a first prototype to explore the possibilities could be built, there is plenty of room for expansion:

- **Replication** - The replication feature is not fully implemented yet, but could be necessary for a fully functional service.

- **Large Data** - Before storing any data, the data itself needs to be held in the application memory for pre-processing. This can result in a huge memory footprint if the data is very large (e.g. big key-value pairs). An improvement could be to start pre-processing during receipt of messages, so the whole data flow then would look more like data streaming: While receiving data, the part already received is already pre-processed on-the-fly and stored to the data store.

- **More Testing** - The algorithms work well, but generally the implementation could be made more stable. In particular, JAVA exceptions are in most cases only caught but not really handled. There is room for improvement to optimize unexpected events (e.g. network connection loss, transmission errors). Currently all unrecognized messages are silently ignored.

- **GUI** - Administration and monitoring of the system could benefit from a graphical interface. Currently, the state of a component can only be checked by looking at the general log file.

- **Security Concerns** - The current version of ERT does not handle any security issues, so malicious systems could interfere with ERT and possibly disable the whole service.

- **Real World Use Case** - A real application that utilizes ERT should be built, to see if this low-level approach is useful.

# List of Figures

# References

[1] L. Allison. Tries, 2004. URL `http://www.csse.monash.edu.au/~lloyd/tildeAlgDS/Tree/Trie/`.

[2] Apple-Computer. Rendezvous, 2003. URL `http://www.apple.com/macosx/features/rendezvous`.

[3] Apple-Computer. Rendezvous Technology Brief, 2003. URL `http://images.apple.com/macosx/pdf/Panther_Rendezvous_TB_10232003.pdf`.

[4] Apple-Computer. Spotlight Technologies, 2004. URL `http://www.apple.com/macosx/tiger/spotlighttech.html`.

[5] Apple-Computer. Xsan for traditional IT services, 2004. URL `http://www.apple.com/xsan/datacenter.html`.

[6] ASN.1 Consortium. ASN.1 Consortium, 2004. URL `http://www.asn1.org`.

[7] J. Aspnes and G. Shah. Skip Graphs. In *Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 384–393, Baltimore, MD, USA, 12–14˜ 2003. URL `http://www.cs.yale.edu/homes/shah/html/pubs/skip-graphs.html`.

[8] H. Balakrishnan and Others. IRIS: Infrastructure for Resilient Internet Systems, 2004. URL `http://iris.lcs.mit.edu`.

[9] S. Brin and L. Page. The Anatomy of a Search Engine, 2004. URL `http://www-db.stanford.edu/~backrub/google.html`.

[10] R. Card, T. Ts'o, and S. Tweedie. Design and Implementation of the Second Extended Filesystem, 2004. URL `http://web.mit.edu/tytso/www/linux/ext2intro.html`.

[11] S. D. S. Center. Distributed Object Computation Testbed (DOCT), 2004. URL `http://www.sdsc.edu/DOCT/Publications/f4-1/f4-1.html`.

[12] CERN. CERN - European Organization for Nuclear Research, 2004. URL `http://www.cern.ch`.

[13] A. Crainiceanu, P. Linga, J. Gehrke, and J. Shanmugasundaram. Querying Peer-to-Peer Networks Using P-Trees, 2004. URL `http://techreports.library.cornell.edu:8081/Dienst/UI/1.0/Display/cul.cis/TR2004-1926`.

[14] A. Crainiceanu, P. Linga, A. Machanavajjhala, J. Gehrke, and J. Shanmugasundaram. P-Ring: An Index Structure for Peer-to-Peer Systems, 2004. URL `http://techreports.library.cornell.edu:8081/Dienst/UI/1.0/Display/cul.cis/TR2004-1946`.

[15] A. Crainiceanu, P. Linga, A. Machanavajjhala, J. Gehrke, and J. Shanmugasundaram. A Storage and Indexing Framework for P2P Systems. Technical report, Cornell University, 2004. URL `http://www.cs.cornell.edu/johannes/papers/2004/www2004-indexingArchitecturePoster.pdf`.

[16] A. Cumming. A Gentle Introduction to SQL, 2003. URL `http://sqlzoo.net/`.

[17] F. Dabek, B. Zhao, P. Druschel, J. Kubiatowicz, and I. Stoica. Towards a Common API for Structured P2P Overlays. In *Proc. of IPTPS*, pages 33–44, Berkeley, CA, Feb 2003. URL `http://www.cs.berkeley.edu/~ravenben/publications/abstracts/apis.html`.

[18] DAMA. DAMA International - Data Facts You Can Use, 2004. URL `http://www.dama.org/public/pages/index.cfm?pageid=461`.

[19] Data & Object Factory. Design Patterns: Proxy, 2004. URL `http://www.dofactory.com/Patterns/PatternProxy.aspx`.

[20] Data & Object Factory. Design Patterns: Singleton, 2004. URL `http://www.dofactory.com/Patterns/PatternSingleton.aspx`.

[21] L. Dongman. Distributed Systems (ICE 601) - Part 1, 2004. URL `http://cds.icu.ac.kr/course/ice601/2004/lecture/intro1.pdf`.

[22] L. Dongman. Distributed Systems (ICE 601) - Part 2, 2004. URL `http://cds.icu.ac.kr/course/ice601/2004/lecture/intro2.pdf`.

[23] P. Druschel and A. Rowstron. Pastry - A substrate for peer-to-peer applications, 2004. URL `http://research.microsoft.com/~antr/Pastry`.

[24] J. Frankel. Gnutella File Sharing, 2004. URL `http://www.gnutella.com`.

[25] E. Fredkin. Trie memory. *Commun. ACM*, 3(9):490–499, 1960. ISSN 0001-0782.

[26] E. Gamma, R. Helm, R. Johnson, and J. Vlissides. *Design Patterns*. Addison-Wesley Professional; 1st edition (January 15, 1995), 1995. ISBN 0201633612.

[27] Getty Research Institute. Introduction to Metadata, 2004. URL `http://www.getty.edu/research/conducting_research/standards/intrometadata/`.

[28] Google. Google Desktop Search, 2004. URL `http://desktop.google.com/`.

[29] Google. Google Print, 2004. URL `http://print.google.com/`.

[30] Google. Google Web Search Features - Cache, 2004. URL `http://www.google.com/help/features.html#cached`.

[31] Hewlett-Packard. Cluster Hardware Configuration, 2002. URL `http://h30097.www3.hp.com/docs/base_doc/DOCUMENTATION/V51B_ACRO_DUX/ARHGWETE.PDF`.

[32] M. Hitchens. Distributed Systems (ITEC 801) - Introduction, 2003. URL `http://www.comp.mq.edu.au/units/itec801/lectures/01-intro-ovs.pdf`.

[33] P. Howard. Database Performance IBM, Oracle and Microsoft an evaluation. Technical report, Bloor Research, 2003. URL `ftp://ftp.software.ibm.com/software/data/pubs/papers/bloor0104.pdf`.

[34] M. Huras. Matt Huras and Version 8 Enhancements for DB2 on Linux, UNIX, and Windows, 2002. URL `http://www-106.ibm.com/developerworks/db2/library/techarticle/0207huras/0207huras.html`.

[35] IETF. Zero Configuration Networking (Zeroconf), 1999. URL `http://www.zeroconf.org`.

[36] O. S. Initiative. Open Source Initiative OSI - The BSD License:Licensing, 2004. URL `http://www.opensource.org/licenses/bsd-license.php`.

[37] G. B. Jacqueline Bloemen. A technical comparison. Technical report, Fourth Millennium Technology Inc. USA, 2003. URL `ftp://ftp.software.ibm.com/software/data/highlights/dbmscomparison.pdf`.

[38] A. K. Jane Wright. Compaq TruCluster Server Unix-Based Clustering Solution. Technical report, Gartner Research, 2002. URL `http://www.totaltec.com/Gartner-Cluster-Solutions.pdf`.

[39] T. Johnson and A. Colbrook. A Distributed Data-Balanced Dictionary Based on the B-Link Tree. Technical Report MIT/LCS/TR-530, 1992. URL `http://citeseer.ist.psu.edu/johnson92distributed.html`.

[40] F. Kaashoek, R. Morris, F. Dabek, I. Stoica, E. Brunskill, D. Karger, R. Cox, and A. Muthitacharoen. The Chord Project, 2004. URL `http://www.pdos.lcs.mit.edu/chord`.

[41] B. Karp, S. Ratnasamy, S. Rhea, and S. Shenker. Spurring Adoption of DHTs with OpenHash, a Public DHT Service, 2004. URL `http://iptps04.cs.ucsd.edu/papers/karp-openhash.pdf`.

[42] J. Kubiatowicz. The OceanStore Project, 2004. URL `http://oceanstore.cs.berkeley.edu`.

[43] I. Kuz. Distributed Systems (COMP 9243) - Introduction, 2004. URL `http://www.cse.unsw.edu.au/~s9243/lectures/intro-slides.pdf`.

[44] I. Kuz. Distributed Systems (COMP 9243) - Notes, 2004. URL `http://www.cse.unsw.edu.au/~s9243/lectures/intro-notes.pdf`.

[45] J. Larmouth, D. Steedman, and J. E. White. ASN.1 Information Site, 2004. URL `http://asn1.elibel.tm.fr/en`.

[46] T. Lee. Netty 2 - Overview, 2004. URL `http://gleamynode.net/dev/projects/netty2/`.

[47] P. L. Lehman and s. Bing Yao. Efficient locking for concurrent operations on B-trees. *ACM Trans. Database Syst.*, 6(4):650–670, 1981. ISSN 0362-5915. URL `http://www.it.iitb.ac.in/~it603/resources/papers/lehman.pdf`.

[48] A. Lisitsa. Distributed System Requirements, 2004. URL `http://www.csc.liv.ac.uk/~alexei/COMP514/COMP514-distributed-systems-04.pdf`.

[49] Microsoft. The Boxwood Project, 2004. URL `http://research.microsoft.com/research/sv/Boxwood/`.

[50] Microsoft. Data Access and Storage Developer Center: Building WinFS Solutions, 2004. URL `http://msdn.microsoft.com/data/winfs/`.

[51] M. Miley. The Grid - bringing computing power to the masses. Oracle Corporation, 2003. URL `http://otn.oracle.com/tech/grid/collateral/Grid_Feature.pdf`.

[52] E. Miller, R. Swick, D. Brickley, B. McBride, J. Hendler, G. Schreiber, and D. Connolly. W3C Semantic Web, 2004. URL `http://www.w3.org/2001/sw/`.

[53] MIMI. Access time - (Computing): Definition, 2004. URL `http://en.mimi.hu/computing/access_time.html`.

[54] NAMESYS. Reiser File System, 2004. URL `http://www.namesys.com/`.

[55] S. Ogura and T. Miura. Paging B-Trees for Distributed Environment. *Workshop on Distributed Data and Structure (WDAS)*, 2002. URL `http://www.ieice.org/iss/de/DEWS/proc/2002/presentations/B2-6.pdf`.

[56] Oracle-Corporation. Tutorial: Grid Computing Demo, 2004. URL `http://otn.oracle.com/tech/grid/index.html`.

[57] V. Pande. Folding@Home Distributed Computing, 2004. URL `http://folding.stanford.edu/`.

[58] W. Pugh. Skip lists: a probabilistic alternative to balanced trees. *Commun. ACM*, 33(6):668–676, 1990. ISSN 0001-0782. URL `ftp://ftp.cs.umd.edu/pub/skipLists/skiplists.pdf`.

[59] S. Rhea. The Bamboo Distributed Hash Table, 2004. URL `http://bamboo-dht.org`.

[60] L. Richter. Verteilte Systeme - Backup Services, 1999. URL `http://www.ifi.unizh.ch/richter/Classes/VS_WS00/V-5-BackupServices.pdf`.

[61] L. Richter. Verteilte Systeme - Einführung, 1999. URL `http://www.ifi.unizh.ch/richter/Classes/VS_WS00/V-1-Einf.pdf`.

[62] L. Richter. Verteilte Systeme - Konzepte und Mechanismen, 1999. URL `http://www.ifi.unizh.ch/richter/Classes/VS_WS00/V-3-Konz&Mech.pdf`.

[63] Y. Sagiv. Concurrent operations on B-trees with overtaking. In *Proceedings of the fourth ACM SIGACT-SIGMOD symposium on Principles of database systems*, pages 28–37. ACM Press, 1985. ISBN 0-89791-153-9.

[64] C. Schmidt and M. Parashar. Flexible Information Discovery in Decentralized Distributed Systems. In *12th IEEE International Symposium on High Performance Distributed Computing*, 2003. URL `http://www.caip.rutgers.edu/~cristins/schmidtc_discovery.pdf`.

[65] Sleepycat-Software. Berkeley DB, 2004. URL `http://www.sleepycat.com/products/db.shtml`.

[66] I. Stoica, R. Morris, D. Karger, M. F. Kaashoek, and H. Balakrishnan. Chord: A Scalable Peer-to-peer Lookup Service for Internet Applications. In *Proceedings of the ACM SIGCOMM '01 Conference*, San Diego, California, August 2001. URL `http://www.pdos.lcs.mit.edu/papers/chord:sigcomm01`.

[67] Sun Microsystems. Java Remote Method Invocation (Java RMI), 2004. URL `http://java.sun.com/products/jdk/rmi/`.

[68] Sun Microsystems. New I/O APIs, 2004. URL `http://java.sun.com/j2se/1.4/nio/`.

[69] A. S. Tanenbaum and M. van Steen. *Distributed systems : principles and paradigms*. Prentice Hall Pearson Education International, Upper Saddle River, NJ, international ed. edition, 2002. ISBN 0-13-1217860(pbk.-).

[70] UC Berkeley. How Much Information, 2004. URL `http://www.sims.berkeley.edu/research/projects/how-much-info-2003/execsum.htm`.

[71] A. van Hoff. JmDNS home page, 2004. URL `http://jmdns.sourceforge.net/`.

[72] M. Welsh, D. E. Culler, and E. A. Brewer. SEDA: An Architecture for Well-Conditioned, Scalable Internet Services. In *Symposium on Operating Systems Principles*, pages 230–243, 2001. URL `citeseer.ist.psu.edu/welsh01seda.html`.

[73] Wikipedia. AVL tree - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/AVL_tree`.

[74] Wikipedia. B-tree - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/B-tree`.

[75] Wikipedia. Backward compatibility - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Backward_compatibility`.

[76] Wikipedia. Binary search tree - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Binary_search_tree`.

[77] Wikipedia. Binary tree - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Binary_tree`.

[78] Wikipedia. Database - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Database`.

[79] Wikipedia. Forward compatibility - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Forward_compatibility`.

[80] Wikipedia. HFS Plus - Wikipedia, 2004. URL `http://en.wikipedia.org/wiki/HFS_Plus`.

[81] Wikipedia. Mainframe computer - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Mainframe`.

[82] Wikipedia. Metadata (computing) - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Metadata_%28computing%29`.

[83] Wikipedia. Network File System - Wikipedia, 2004. URL `http://en.wikipedia.org/wiki/NFS`.

[84] Wikipedia. NTFS - Wikipedia, 2004. URL `http://en.wikipedia.org/wiki/NTFS`.

[85] Wikipedia. Red-black tree - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Red-black_tree`.

[86] Wikipedia. Relational model - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Relational_model`.

[87] Wikipedia. Self-balancing binary search tree - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Self-balancing_binary_search_tree`.

[88] Wikipedia. Splay tree - Wikipedia, the free encyclopedia, 2004. URL `http://en.wikipedia.org/wiki/Splay_tree`.

[89] P. Yalagandula. Solving Range Queries in a Distributed System. Technical report, University of Texas at Austin, 2003. URL `http://www.cs.utexas.edu/users/browne/cs395f2003/projects/YalagandalaReport.pdf`.

[90] B. Y. Zhao, A. Joseph, and J. Kubiatowicz. Tapestry - Infrastructure for Fault-resilient, Decentralized Location and Routing, 2004. URL `http://www.cs.berkeley.edu/~ravenben/tapestry`.

# A Storage Subsystems

Using a storage subsystem has several advantages over directly accessing the disk:

- More than one node can directly access data - otherwise, only one node can mount the physical hard-drive.

- Because of an HAL (Hardware Abstraction Layer), the actual implementation of a storage subsystem can abstract several hardware limitations normally found on direct disk-based access:

  - Redundancy: If one disk fails, another can take over (i.e. RAID)

  - Adding/Removing/Reallocating disk space on the fly

  - Enhanced Caching mechanisms

- Storage is in one place; this makes it easier to administrate.

## A.1 TruCluster Shared Storage

**Hewlett-Packard Corporation** - TruCluster is the cluster offering from HP (former Compaq) and also contains a shared storage mechanism.



Figure 54: Shared Storage Overview (taken from [31])

Here one cluster member takes the lead and mounts the cluster-wide shared-storage disk. All other cluster members need to send all read/write requests to this

leading member over the cluster-interconnect. If the leader is down, there is a sophisticated leader-elector process (not shown here) and the new leader then mounts the cluster-wide file system. This mechanism will not scale for a larger number of nodes. In general, truCluster is limited to 8 nodes in the current version.

## A.2   Xsan

**Apple Computer Corporation, Advanced Digital Information Corporation** - As a newcomer to this market, they have a simple offering which helps consolidate storage to a common place. It uses technology from ADIC (Advanced Digital Information Corporation).

This approach also has limitations and will not scale indefinitely. The number of Storage Servers is limited (64 nodes) and also the concurrent access to the physical disks is limited. The overall performance could be higher than the truCluster approach, because it uses a dedicated storage access medium (Fibre Channel) used only for transmitting storage related data, where in truCluster, in addition to storage-related data, general cluster information data is also transmitted.



Figure 55: Xsan Overview (taken from [5])

1. **Storage Pool** - The actual disk drives.

2. **Metadata Controller** - Manages the Storage Servers (3) and directs them to the actual files in the Storage Pool.

3. **Storage Servers** - Are connected to the Storage Pool using Fibre Channel and use the Metadata Controller (2) for concurrent file access.

4. **Data Consumers** - Clients who want to access the data.

5. **Fibre Channel Switch** - Dedicated switch for Storage Pool (1).

6. **Backup Metadata Controller** - Provides high-availability insurance.

# B   Content of CD-ROM

The attached CD-ROM contains the following data:

- **Abstract** - A file containing the English abstract of this thesis.

- **Zusfsg** - A file containing the German abstract of this thesis.

- **code** - A folder containing the most recent version of the ERT software, including instructions on how to install and run the software.

- **benchmark** - A folder containing the scripts for automated testing, test protocol and a system profile of one of the test machines.

- **misc** - A folder containing additional files that have been produced during writing this thesis. This includes copies of PDF documents found during research and intermediate papers handed in during writing the thesis.

- **repository** - A folder containing the software repository for the software and the thesis paper itself. During the writing of the thesis, all intermediate steps were recorded into this repository. This allows retrieval of the thesis and source code at any given date. The repository is stored as a Subversion archive (see `http://subversion.tigris.org`) in two different formats. One as native format and the other as a dump file, to import the repository into another existing Subversion repository.

- **thesis** - A folder containing the final version of the thesis as PDF, LaTeX, LyX and plain text document.

# C   Source Code

**File 1:** base.Dispatcher

```
 1 /*
 2  * Created on Aug 13, 2004 11:44:01 AM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:   $Author: michi $
 7  * Revision:     $Rev: 138 $
 8  * Location:     $URL: Dispatcher.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.base;
17
18 import java.util.Hashtable;
19
20 import org.apache.log4j.Logger;
21
22 import ch.nix.ert.message.DataItem;
23 import ch.nix.ert.util.Queue;
24
25 import EDU.oswego.cs.dl.util.concurrent.LinkedQueue;
26
27 /**
28  * TODO Edit this text for a better description
29  *
30  * This is a local dispatcher, he is responsible for all local requests. He does
31  * not decide if the request is executed locally or remotely, he will ALWAYS
32  * execute it remotly. It's the responsibility of each single service to relay
33  * it to another node... <br>
34  * If you don't know, if the request is for local or remote execution, you can
35  * place it in Dispatcher, but he will forward it to KeyPartitionManager who
36  * will then decide if executed locally or remotely (according their KeyGroupID)
37  *
38  * @author michi
39  */
40 public class Dispatcher implements Runnable {
41
42     // logger:
43     static Logger logger = Logger.getLogger(Dispatcher.class);
44
45     private LinkedQueue dispatcherQueue = new LinkedQueue();
46
47     private Hashtable dispatcherServices = new Hashtable();
48
49     private Hashtable serviceNames = new Hashtable();
50
51     private Object queueLock = new Object();
52
53     /**
54      * A handle to the unique Dispatcher instance.
55      */
56     static private Dispatcher instance = null;
57
58     private Dispatcher() {
59     }
60
61     static public synchronized Dispatcher getInstance() {
62         if (null == instance) {
63             instance = new Dispatcher();
64             logger.info("Dispatcher started");
65         }
66         return instance;
67     }
68
69     /**
```

```
 70        * Only data which contain the "command" key are accepted, all other will be
 71        * silently ignored! TODO: perhpas we could also throw an exception?
 72        *
 73        * @param data
 74        */
 75       public void addWork(DataItem data) {
 76           synchronized (queueLock) {
 77               //check data, and only add it to queue, if it contains a operation
 78               // (like "put")
 79               if (data.getOperation() != null) {
 80                   try {
 81                       //System.out.println("DispatcherWork: "+data);
 82                       dispatcherQueue.put(data);
 83                   } catch (InterruptedException e) {
 84                       // TODO Auto-generated catch block
 85                       e.printStackTrace();
 86                   }
 87               }
 88           }
 89       }
 90
 91       /**
 92        * There is only ONE service allowed for each type... if an service is
 93        * already registered and a new one with the same keyword, subtype is added,
 94        * an exception is thrown TODO: implment this EXCEPTION
 95        *
 96        * @param command
 97        *            The events you are interested in (e.g. "put", "get",
 98        *            "delete","ping" <br>
 99        *            Some events have subtypes, subtypes can be: request, response
100        *            (e.g. a "get" request can be of type "request", this means it
101        *            has not been processed yet or of type "response" this means
102        *            this is the response of the first request <br>
103        *            <i>request </i>starting of a transaction <br>
104        *            <i>response </i>ending of a transaction
105        * @param queue
106        *            The queue where incomming request are stored for this command.
107        */
108       public void registerService(String operation, Queue queue) {
109           synchronized (dispatcherServices) {
110               operation = operation.toLowerCase();
111               dispatcherServices.put(operation, queue);
112
113           }
114       }
115
116       public void registerQueue(Queue queue) {
117           serviceNames.put(queue.getName(), queue);
118       }
119
120       public void unregisterService(String operation) {
121           synchronized (dispatcherServices) {
122               operation = operation.toLowerCase();
123               dispatcherServices.remove(operation);
124           }
125       }
126
127       /**
128        * Start the dispatcher and ONLY ONE thread meight be inside it!
129        */
130       public void run() {
131           // First set the name for this Thread, helps debugging :-)
132           Thread.currentThread().setName("ERT.Dispatcher");
133           TransactionMonitor transactionMonitor = new TransactionMonitor();
134           DataItem dataToDispatch = null;
135           String operation;
136           String transactionType;
137           String routingHint;
138           int transactionSplit;
139           while (true) {
140               Queue operationQueue = null;
141               // wait for something in dispatcherQueue
142               try {
143                   dataToDispatch = (DataItem) dispatcherQueue.take();
144               } catch (InterruptedException e) {
145                   // TODO Auto-generated catch block
146                   e.printStackTrace();
147               }
148               // now dispatch request
149               if (dataToDispatch != null) {
150
151                   // transactionType = dataToDispatch.getTransactionType();
152                   routingHint = dataToDispatch.getRouting();
153
154                   // reset the transactionSplit
```

```
155                    transactionSplit = 0;
156
157                // now find the correct ServiceQueue where to put the value
158                if (routingHint == null || routingHint.equals("local")) {
159                    transactionType = dataToDispatch.getTransactionType();
160                    if (transactionType != null
161                            && transactionType.equals("end")) {
162                        // we send the response right to the place where it
163                        // originated
164                        operationQueue = (Queue) serviceNames
165                                .get(dataToDispatch.getTransactionServiceId());
166
167                        // now check, if we this transaction is part to a
168                        // splitted transaction
169                        transactionSplit = dataToDispatch.getTransactionSplit();
170                        if (transactionSplit > 0) {
171                            // yes, it is a result from a splited transaction
172                            // pass it on to the TransactionMonitor whoch holds
173                            // references to all open splited transactions
174                            dataToDispatch = transactionMonitor
175                                    .put(dataToDispatch);
176                            if (dataToDispatch == null) {
177                                operationQueue = null;
178                            }
179                        }
180                    } else {
181                        // as I already checked at insertion time, this data has
182                        // an operation defined!
183                        operation = dataToDispatch.getOperation();
184                        operationQueue = (Queue) dispatcherServices
185                                .get(operation);
186                    }
187                } else {
188                    // default means mostly the item is sent to the
189                    // NodeCommunication layer to sending out to another
190                    // node
191                    operationQueue = (Queue) dispatcherServices.get("default");
192                }
193
194                if (operationQueue != null) {
195                    try {
196                        //System.out.println("Data2Dispatch: "+dataToDispatch+"
197                        // is response? "+dataToDispatch.isRespose());
198                        operationQueue.put(dataToDispatch);
199 //                        logger.debug("Work: " + dataToDispatch
200 //                                + " put on queue: " + operationQueue.getName());
201                    } catch (InterruptedException e) {
202                        // TODO Auto-generated catch block
203                        e.printStackTrace();
204                    }
205                } else if (transactionSplit == 0) {
206                    logger.warn("No Queue found for dataItem: "
207                            + dataToDispatch);
208                }
209            }
210        }
211    }
212 }
```

---

**File 2:** base.RootLogger

```
1 /*
2  * Created on Sep 7, 2004 11:05:37 AM
3  *
4  * Project:     ExtendendRootTree
5  * Last Change: $Date: 2004-11-11 18:55:48 +0100 (Thu, 11 Nov 2004) $
6  * Changed by:  $Author: michi $
7  * Revision:    $Rev: 169 $
8  * Location:    $URL: RootLogger.java $
9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.base;
17
18 import org.apache.log4j.BasicConfigurator;
19 import org.apache.log4j.DailyRollingFileAppender;
20 import org.apache.log4j.Level;
21 import org.apache.log4j.Logger;
```

```
22 import org.apache.log4j.PatternLayout;
23
24 /**
25  * TODO Edit this text for a better description
26  *
27  * @author michi
28  */
29 public class RootLogger {
30
31     static Logger logger = Logger.getRootLogger();
32
33     /**
34      *
35      */
36     public RootLogger() {
37         // TTCCLayout layout = new TTCCLayout();
38         PatternLayout layout = new PatternLayout(
39                 "%d{HH:mm:ss,SSS} [%15.15t] %-5p %30.30c.%-3L %x - %m%n");
40         DailyRollingFileAppender appender = null;
41         try {
42             appender = new DailyRollingFileAppender(layout,
43                     "/Users/tutor/ert/log/ert.log", "'.'yyyy-MM-dd");
44         } catch (Exception e) {
45         }
46
47         logger.setLevel(Level.INFO);
48         BasicConfigurator.configure(appender);
49         logger.info("-------- ERT started --------");
50     }
51 }
```

---

**File 3:** base.ServiceBase

```
 1 /*
 2  * Created on Sep 8, 2004 7:43:33 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 138 $
 8  * Location:    $URL: ServiceBase.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.base;
17
18 import org.apache.log4j.Logger;
19
20 import ch.nix.ert.message.DataItem;
21 import ch.nix.ert.message.Node;
22 import ch.nix.ert.util.Queue;
23
24 /**
25  * TODO Edit this text for a better description
26  *
27  * @author michi
28  */
29 public abstract class ServiceBase implements ServiceInterface {
30
31     // logger:
32     static Logger logger = Logger.getLogger(ServiceBase.class);
33
34     private Dispatcher dispatcher;
35
36     private Queue serviceQueue;
37
38     private int transactionId = 0;
39
40     private Node currentNode;
41
42     private String serviceName;
43
44     private int splitSize = 0;
45
46     private String storedId;
47
48     private Node storedNodeId;
49
```

```
50      private String storedServiceId;
51
52      private String storedOperation;
53
54      private String storedTransactionType;
55
56      private String storedTransactionOperation;
57
58      /**
59       *
60       */
61      protected ServiceBase() {
62          //this.getClass().getName().
63          dispatcher = ServiceManager.getInstance().getDispatcher();
64          currentNode = new Node();
65          serviceName = this.getClass().getName();
66          serviceQueue = new Queue(getQueueName());
67          // register service Queue, so dispatcher knows where to send the
68          // responses
69          dispatcher.registerQueue(serviceQueue);
70      }
71
72      protected void registerService(String service) {
73          dispatcher.registerService(service, serviceQueue);
74      }
75
76      public void addWork(DataItem data) {
77          dispatcher.addWork(data);
78      }
79
80      //    public void transactionBeginStatefull(DataItem data) {
81      //        data.setTransactionId(getTransactionId());
82      //        transactionBeginGeneral(data);
83      //    }
84      //
85      //    public void transactionBegin(DataItem data) {
86      //        //data.setTransactionId(data.getOperation());
87      //        transactionBeginGeneral(data);
88      //    }
89
90      public void transactionBegin(DataItem data) {
91          data.setTransactionNodeId(getNodeId());
92          data.setTransactionServiceId(getServiceId());
93          data.setTransactionId(getTransactionId());
94          data.setTransactionOperation(data.getOperation());
95          data.setTransactionType("begin");
96          dispatcher.addWork(data);
97      }
98
99      public void transactionContinue(DataItem data) {
100         // only fill in data, if incoming request really was a transaction,
101         // otherwise just send back the answer
102         if (storedTransactionType != null) {
103             data.setTransactionType("continue");
104             data.setTransactionSplit(splitSize);
105             data.setTransactionNodeId(storedNodeId);
106             data.setTransactionServiceId(storedServiceId);
107             data.setTransactionId(storedId);
108             //data.setTransactionOperation(storedOperation);
109         }
110         dispatcher.addWork(data);
111     }
112
113     public void transactionSplit(DataItem data, int newSplitSize) {
114         // only fill in data, if incoming request really was a transaction,
115         // otherwise just send back the answer
116         if (storedTransactionType != null) {
117             data.setTransactionType("split");
118             data.setTransactionSplit(splitSize+newSplitSize);
119             data.setTransactionNodeId(storedNodeId);
120             data.setTransactionServiceId(storedServiceId);
121             data.setTransactionId(storedId);
122             //data.setTransactionOperation(storedOperation);
123         }
124         dispatcher.addWork(data);
125     }
126
127     public void transactionEnd(DataItem data) {
128         // only fill in data, if incoming request really was a transaction,
129         // otherwise we don't send back anything, as a transaction does not make
130         // any sense...
131         if (storedTransactionType != null) {
132             data.setTransactionType("end");
133             data.setTransactionSplit(splitSize);
134             data.setTransactionNodeId(storedNodeId);
```

```
135              data.setTransactionServiceId(storedServiceId);
136              data.setTransactionId(storedId);
137              data.setOperation(storedTransactionOperation);
138
139              if (data.isTransactionLocal()) {
140                  data.setRouting("local");
141              } else {
142                  data.setDirectDestination(data.getTransactionNodeId());
143                  data.setRouting("direct");
144              }
145              dispatcher.addWork(data);
146          }
147      }
148
149      public Queue getServiceQueue() {
150          return serviceQueue;
151      }
152
153      public DataItem getWork() {
154          try {
155              DataItem data = (DataItem) serviceQueue.take();
156              //String transactionType = data.getTransactionType();
157              //logger.warn("getWork: "+data);
158              // if (transactionType != null && transactionType.equals("end")) {
159              //logger.warn("if is executed");
160              //    data.setOperation(data.getOperation() + "-end");
161              //}
162              storedOperation = data.getOperation();
163              splitSize = data.getTransactionSplit();
164              storedId = data.getTransactionId();
165              storedServiceId = data.getTransactionServiceId();
166              storedNodeId = data.getTransactionNodeId();
167              storedTransactionType = data.getTransactionType();
168              storedTransactionOperation = data.getTransactionOperation();
169
170              return data;
171          } catch (InterruptedException e) {
172              // TODO Auto-generated catch block
173              e.printStackTrace();
174          }
175          return null;
176      }
177
178      public DataItem getWork(long timeout) {
179          try {
180              return (DataItem) serviceQueue.poll(timeout);
181          } catch (InterruptedException e) {
182              // TODO Auto-generated catch block
183              e.printStackTrace();
184          }
185          return null;
186      }
187
188      //TODO do a nice regex search...
189      private String getQueueName() {
190          return this.getClass().getName();
191          //System.out.println(this.getClass().getName().split("\\..*Service$"));
192      }
193
194      private Node getNodeId() {
195          return currentNode;
196      }
197
198      private String getServiceId() {
199          return this.getClass().getName();
200      }
201
202      private String getTransactionId() {
203          return Integer.toString(transactionId++);
204      }
205
206 }
```

---

**File 4:** base.ServiceInterface

---

```
1 /*
2  * Created on Aug 13, 2004 3:03:03 PM
3  *
4  * Project:     ExtendendRootTree
5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
6  * Changed by:  $Author: michi $
7  * Revision:    $Rev: 138 $
```

```
 8  * Location:     $URL: ServiceInterface.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.base;
17
18 /**
19  * TODO Edit this text for a better description
20  *
21  * @author michi
22  */
23 public interface ServiceInterface extends Runnable {
24
25      /**
26       * Registers the Service at the Dispatcher
27       */
28      public void registerAtDispatcher();
29
30      // public Queue getServiceQueue();
31
32 }
```

## File 5: base.ServiceManager

```
 1 /*
 2  * Created on Aug 6, 2004 1:37:14 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:43:26 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 139 $
 8  * Location:    $URL: ServiceManager.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.base;
17
18 import java.lang.reflect.InvocationTargetException;
19 import java.lang.reflect.Method;
20 import java.util.ArrayList;
21
22 import org.apache.log4j.Logger;
23
24 /**
25  * TODO Edit this text for a better description <br>
26  *
27  * Follows the Singelton-Design Pattern after Rodney Waldhoff.
28  * http://radio.weblogs.com/0122027/stories/2003/10/20/implementingTheSingletonPatternInJava.html
29  *
30  * @author michi
31  */
32 public class ServiceManager {
33
34      //logger:
35      static Logger logger = Logger.getLogger(ServiceManager.class);
36
37      /**
38       * A handle to the unique ServiceManager instance.
39       */
40      static private ServiceManager instance = null;
41
42      static String[] serviceNameList = { "ch.nix.ert.main.MainService",
43              "ch.nix.ert.callback.CallbackService",
44              "ch.nix.ert.communication.NodeCommunicationService",
45              "ch.nix.ert.datastore.DataStoreService",
46              "ch.nix.ert.keypartition.KeyPartitionService",
47              "ch.nix.ert.user.UserService" };
48
49      static ArrayList serviceList = new ArrayList();
50
51      static private Thread dispatcherThread;
52
53      static private Dispatcher dispatcher;
54
```

```
55      //private RootLogger ertLogger;
56
57      /**
58       * The constructor could be made private to prevent others from
59       * instantiating this class. But this would also make it impossible to
60       * create instances of ServiceManager subclasses.
61       */
62      private ServiceManager() {
63          // TODO temporarely create logger here, afterwards it needs to be put
64          // into startup class
65          // ertLogger = new RootLogger();
66      }
67
68      /**
69       * @return The unique instance of this class.
70       */
71      static public synchronized ServiceManager getInstance() {
72          if (null == instance) {
73              instance = new ServiceManager();
74              dispatcher = Dispatcher.getInstance();
75              dispatcherThread = new Thread(dispatcher);
76              String serviceName;
77              for (int i = 0; i < serviceNameList.length; i++) {
78                  serviceName = serviceNameList[i];
79                  ServiceInterface service;
80                  Thread serviceThread;
81                  Method method;
82
83                  try {
84                      method = (Class.forName(serviceName).getMethod(
85                          "getInstance", null));
86                      service = (ServiceInterface) method.invoke(null, null);
87                      service.registerAtDispatcher();
88                      serviceThread = new Thread(service);
89                      serviceList.add(serviceThread);
90                      serviceThread.start();
91                      logger.debug("Service loaded & started: " + serviceName);
92                  } catch (SecurityException e) {
93                      logger.error("Error creating: " + serviceName + "\n" + e);
94                  } catch (NoSuchMethodException e) {
95                      logger.error("Error creating: " + serviceName + "\n" + e);
96                  } catch (ClassNotFoundException e) {
97                      logger.error("Error creating: " + serviceName + "\n" + e);
98                  } catch (IllegalArgumentException e) {
99                      logger.error("Error creating: " + serviceName + "\n" + e);
100                 } catch (IllegalAccessException e) {
101                     logger.error("Error creating: " + serviceName + "\n" + e);
102                 } catch (InvocationTargetException e) {
103                     logger.error("Error creating: " + serviceName + "\n" + e);
104                 }
105
106             }
107
108             //and lastely:
109             dispatcherThread.start();
110             logger.info("ServiceManager started");
111         }
112         return instance;
113     }
114
115     public Dispatcher getDispatcher() {
116         return dispatcher;
117     }
118
119 }
```

**File 6:** base.TransactionMonitor

```
1 /*
2  * Created on Sep 24, 2004 4:14:16 PM
3  *
4  * Project:      ExtendendRootTree
5  * Last Change: $Date: 2004-10-26 15:43:26 +0200 (Tue, 26 Oct 2004) $
6  * Changed by:  $Author: michi $
7  * Revision:    $Rev: 139 $
8  * Location:    $URL: TransactionMonitor.java $
9  *
10 * This source-code is part of the diploma thesis of Michael Kussmaul
11 * written at the University of Zurich, Switzerland
12 * Department of Information Technology (www.ifi.unizh.ch)
13 * For copyright information please contact kussmaul@nix.ch
14 *
```

```
15  */
16 package ch.nix.ert.base;
17
18 import org.apache.log4j.Logger;
19
20 import ch.nix.ert.message.DataItem;
21 import ch.nix.ert.util.LinkedHashMap;
22
23 /**
24  * Stores the temporary results for nested (or split) Transactions. Only if all
25  * parts are arrived, this TransactionMonitor will join the results and return
26  * the transaction
27  * <p>
28  *
29  * In current implementation only boolean operations are allowed to be
30  * split/joined (e.g. those who have ok/Error as result)
31  *
32  * @author michi
33  */
34 public class TransactionMonitor {
35
36     // logger:
37     static Logger logger = Logger.getLogger(TransactionMonitor.class);
38
39     // if more than 10000 open transactions, the oldest one will be discarded...
40     private int bufferSize = 10000;
41
42     private LinkedHashMap transactionIds;
43
44     public TransactionMonitor() {
45         transactionIds = new LinkedHashMap();
46         transactionIds.setMaxSize(bufferSize);
47     }
48
49     /**
50      * @param data
51      * @return The final transaction DataItem, if we have now all parts
52      *         together. It contains also an variable called result, which
53      *         indicates, if operation was a success or not. If it was only a
54      *         part of a transaction, it will return null
55      */
56     public DataItem put(DataItem dataItem) {
57         // logger.error("Monitor: got it: "+dataItem);
58         // first check the result type of this part, if it is already false, we
59         // can already return this dataItem, as the whole transaction is wrong
60         if (dataItem.getResult() == false) {
61             removeId(dataItem.getTransactionId());
62             return dataItem;
63         }
64
65         String transactionId = dataItem.getTransactionId();
66         addOneToId(transactionId);
67
68         int splitAmount = dataItem.getTransactionSplit();
69         if (splitAmount <= getAmountForId(transactionId)) {
70             // means we have all parts together and can return success :-)
71             return dataItem;
72         } else {
73             return null;
74         }
75
76     }
77
78     private int getAmountForId(String id) {
79         return ((Integer) transactionIds.get(id)).intValue();
80     }
81
82     private void addOneToId(String id) {
83         Integer value = (Integer) transactionIds.get(id);
84         int amount;
85         if (value == null) {
86             amount = 0;
87         } else {
88             amount = value.intValue();
89         }
90         transactionIds.put(id, new Integer(amount + 1));
91     }
92
93     private void removeId(String id) {
94         transactionIds.remove(id);
95     }
96 }
```

**File 7:** callback.CallbackService

```
 1 /*
 2  * Created on Aug 18, 2004 12:55:22 PM
 3  *
 4  * Project:       ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:     $Rev: 138 $
 8  * Location:     $URL: CallbackService.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.callback;
17
18 import java.util.Calendar;
19 import java.util.TreeMap;
20
21 import org.apache.log4j.Logger;
22
23 import ch.nix.ert.base.ServiceBase;
24 import ch.nix.ert.message.DataItem;
25
26 /**
27  * TODO Edit this text for a better description
28  *
29  * @author michi
30  */
31 public class CallbackService extends ServiceBase {
32
33     // logger:
34     static Logger logger = Logger.getLogger(CallbackService.class);
35
36     static private CallbackService instance = null;
37
38     //private Queue serviceQueue = new Queue("CallbackQueue");
39
40     //private Dispatcher dispatcher;
41
42     private TreeMap timeSortedList = new TreeMap();
43
44     private CallbackService() {
45     }
46
47     static public synchronized CallbackService getInstance() {
48         if (null == instance) {
49             instance = new CallbackService();
50             //instance.dispatcher =
51             // ServiceManager.getInstance().getDispatcher();
52             logger.info("CallbackService started");
53         }
54         return instance;
55     }
56
57     /*
58      * (non-Javadoc)
59      *
60      * @see ch.nix.ert.ServiceInterface#registerAtDispatcher()
61      */
62     public void registerAtDispatcher() {
63         // TODO Auto-generated method stub
64         registerService("calllater");
65         registerService("callcontinous");
66     }
67
68     /*
69      * (non-Javadoc)
70      *
71      * @see java.lang.Runnable#run()
72      */
73     public void run() {
74         Thread.currentThread().setName("ERT.CallbackService");
75         // get the first request
76         long currentTime, wait;
77         Long newestTime;
78         DataItem dataItem = null;
79         dataItem = getWork();
80         // System.out.println("CallBackService got: "+dataItem);
81         logger.debug("Register Callback: " + dataItem);
82         this.putInSortedTimeList(dataItem);
```

```
 83          while (true) {
 84              currentTime = Calendar.getInstance().getTimeInMillis();
 85              newestTime = (Long) timeSortedList.firstKey();
 86              wait = newestTime.longValue();
 87              if (newestTime.longValue() < currentTime) {
 88                  dataItem = (DataItem) timeSortedList.remove(newestTime);
 89                  DataItem dataToExecute = dataItem.getNested();
 90                  // fill in transaction-Data from sourounding dataItem:
 91                  if (dataItem.getTransactionType() != null) {
 92                      // request was a transaction, so i Have to swap in all data
 93                      // from the main object into the nested one...
 94                      dataToExecute.setTransactionType("continue");
 95                      dataToExecute.setTransactionNodeId(dataItem
 96                              .getTransactionNodeId());
 97                      dataToExecute.setTransactionServiceId(dataItem
 98                              .getTransactionServiceId());
 99                      dataToExecute.setTransactionId(dataItem.getTransactionId());
100                      dataToExecute.setTransactionOperation(dataToExecute
101                              .getOperation());
102                      transactionContinue(dataToExecute);
103                  } else {
104                      addWork(dataToExecute);
105                  }
106                  logger.debug("Execute Callback: " + dataItem);
107                  //now check, if this is a continous service, so we have to put
108                  // it into queue again...
109                  if (dataItem.getOperation().equals("callcontinous")) {
110                      this.putInSortedTimeList(dataItem);
111                  }
112                  if (timeSortedList.isEmpty()) {
113                      wait = -1;
114                  } else {
115                      wait = ((Long) timeSortedList.firstKey()).longValue();
116                      //System.out.println("Set Waiting to: "+wait);
117                  }
118              }
119              if (wait == -1) {
120                  // means there is no element in the queue anymore, so I have
121                  // to wait indefinetly:
122                  dataItem = getWork();
123                  logger.debug("Register Callback: " + dataItem);
124                  // System.out.println("CallBackService got: "+dataItem);
125              } else {
126                  dataItem = getWork(wait - currentTime);
127                  logger.debug("Register Callback or Timeout: " + dataItem);
128                  // System.out.println("CallBackService got: "+dataItem);
129              }
130              if (dataItem != null) {
131                  this.putInSortedTimeList(dataItem);
132              }
133
134          }
135      }
136
137      private void putInSortedTimeList(DataItem data) {
138          long currentTime = Calendar.getInstance().getTimeInMillis();
139          long delay = data.getDelay();
140          long nextTime = currentTime + delay;
141          //TODO: Performance enhancements todo...
142          //As Java Maps/Sets do not allow duplicated keys, I have to check
143          // first, if this key is already in queue... if so I have to postpone it
144          // a bit... this is not very performant, but this class is anyway not
145          // used that often in the application. If you rely on exact timing, a
146          // external Priority-Queue class needs to be built, which supports
147          // duplicated keys
148          boolean duplicateFound = true;
149          while (duplicateFound) {
150              if (timeSortedList.containsKey(new Long(nextTime))) {
151                  nextTime++;
152              } else {
153                  duplicateFound = false;
154              }
155          }
156          timeSortedList.put(new Long(nextTime), data);
157      }
158
159 }
```

---

**File 8:** communication.NodeCommunication

```
1 /*
2  * Created on Aug 10, 2004 11:26:59 AM
```

```
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 138 $
 8  * Location:    $URL: NodeCommunication.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication;
17
18 import org.apache.log4j.Logger;
19
20 import ch.nix.ert.message.DataItem;
21
22 /**
23  * TODO Edit this text for a better description
24  *
25  * Follows the Singelton-Design Pattern after Rodney Waldhoff.
26  * http://radio.weblogs.com/0122027/stories/2003/10/20/implementingTheSingletonPatternInJava.html
27  *
28  * @author michi
29  */
30 public class NodeCommunication implements NodeCommunicationInterface {
31
32     static Logger logger = Logger.getLogger(NodeCommunication.class);
33
34     static private String plugin = "netty2";
35
36     static private NodeCommunication instance = null;
37
38     private NodeCommunicationInterface delegate = null;
39
40     private NodeCommunication() {
41     }
42
43     static public synchronized NodeCommunication getInstance() {
44         if (null == instance) {
45             instance = new NodeCommunication();
46             plugin = plugin.toLowerCase();
47             try {
48                 instance.delegate = (NodeCommunicationInterface) (Class
49                         .forName("ch.nix.ert.communication." + plugin
50                                 + ".NodeCommunication").newInstance());
51                 logger.info("NodeCommunication loaded plugin: "+plugin);
52                 return instance;
53             } catch (InstantiationException e) {
54                 // TODO Auto-generated catch block
55                 e.printStackTrace();
56             } catch (IllegalAccessException e) {
57                 // TODO Auto-generated catch block
58                 e.printStackTrace();
59             } catch (ClassNotFoundException e) {
60                 // TODO Auto-generated catch block
61                 e.printStackTrace();
62             }
63             return null;
64         }
65         return instance;
66
67     }
68
69     /*
70      * (non-Javadoc)
71      *
72      * @see ch.nix.ert.communication.NodeCommunicationInterface#send(ch.nix.ert.DataItem)
73      */
74     public void send(DataItem dataItem) {
75         // TODO Auto-generated method stub
76         delegate.send(dataItem);
77     }
78 }
```

---

**File 9:** communication.NodeCommunicationInterface

---

```
 1 /*
 2  * Created on Aug 10, 2004 11:27:16 AM
 3  *
```

```
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 138 $
 8  * Location:    $URL: NodeCommunicationInterface.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication;
17
18 import ch.nix.ert.message.DataItem;
19
20 /**
21  * TODO Edit this text for a better description
22  *
23  * @author michi
24  */
25 public interface NodeCommunicationInterface {
26
27      public void send(DataItem data);
28
29 }
```

## File 10: communication.NodeCommunicationService

```
 1 /*
 2  * Created on Aug 17, 2004 11:51:59 AM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 138 $
 8  * Location:    $URL: NodeCommunicationService.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication;
17
18 import org.apache.log4j.Logger;
19
20 import ch.nix.ert.base.ServiceBase;
21 import ch.nix.ert.message.DataItem;
22
23 /**
24  * TODO Edit this text for a better description
25  *
26  * @author michi
27  */
28 public class NodeCommunicationService extends ServiceBase {
29
30     // logger:
31     static Logger logger = Logger.getLogger(NodeCommunicationService.class);
32
33     static private NodeCommunicationService instance = null;
34
35     private NodeCommunication nodeCommunication;
36
37     private NodeDirectory nodeDirectory;
38
39     private NodeCommunicationService() {
40     }
41
42     static public synchronized NodeCommunicationService getInstance() {
43         if (null == instance) {
44             instance = new NodeCommunicationService();
45             instance.nodeCommunication = NodeCommunication.getInstance();
46             instance.nodeDirectory = NodeDirectory.getInstance();
47             logger.info("NodeCommunicationService started");
48         }
49         return instance;
50     }
51
52     /*
53      * (non-Javadoc)
```

```
54        *
55        * @see ch.nix.ert.ServiceInterface#registerAtDispatcher()
56        */
57       public void registerAtDispatcher() {
58           registerService("send");
59           registerService("publish");
60           registerService("unpublish");
61
62           registerService("default");
63       }
64
65       /*
66        * (non-Javadoc)
67        *
68        * @see java.lang.Runnable#run()
69        */
70       public void run() {
71           //       System.out.println("DataStoreService started");
72           Thread.currentThread().setName("ERT.NodeCommunicationService");
73           while (true) {
74               DataItem dataItem = null;
75               dataItem = getWork();
76               // System.out.println("got request: " + dataItem);
77               String operation = dataItem.getOperation();
78               if (operation.equals("send")) {
79                   // now check if dataItem is nested. otherwise send the whole
80                   // datItem as it is.
81                   // This can happen, if we automatically send back the response
82                   // to a request, then it meight not have a nested dataItem, as
83                   // nobody has wrapped up the packet.
84                   DataItem dataToSend = dataItem.getNested();
85                   if (dataToSend == null) {
86                       dataToSend = dataItem;
87                   }
88                   //System.out.println("Comm Service2send: " + dataToSend);
89                   nodeCommunication.send(dataToSend);
90               } else if (operation.equals("publish")) {
91                   nodeDirectory.publish(dataItem);
92               } else if (operation.equals("unpublish")) {
93                   nodeDirectory.unpublish(dataItem);
94               } else {
95
96                   //TODO will be uncommented later!
97                   DataItem dataToSend = dataItem.getNested();
98                   if (dataToSend == null) {
99                       dataToSend = dataItem;
100                  }
101                  // System.out.println("Comm Service: " + dataToSend);
102                  nodeCommunication.send(dataToSend);
103              }
104          }
105      }
106 }
```

---

**File 11:** communication.NodeDirectory

```
 1 /*
 2  * Created on Aug 8, 2004 8:11:34 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 138 $
 8  * Location:    $URL: NodeDirectory.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication;
17
18 import org.apache.log4j.Logger;
19
20 import ch.nix.ert.message.DataItem;
21
22 /**
23  * TODO Edit this text for a better description
24  *
25  * Follows the Singelton-Design Pattern after Rodney Waldhoff.
26  * http://radio.weblogs.com/0122027/stories/2003/10/20/implementingTheSingletonPatternInJava.html
```

```
27  *
28  * In addition this class acts as a proxy to the actual implementation.
29  *
30  * @author michi
31  */
32 public class NodeDirectory implements NodeDirectoryInterface {
33
34     // logger:
35     static Logger logger = Logger.getLogger(NodeDirectory.class);
36
37     static private String plugin = "jmdns";
38
39     static private NodeDirectory instance = null;
40
41     private NodeDirectoryInterface delegate = null;
42
43     private NodeDirectory() {
44     }
45
46     static public synchronized NodeDirectory getInstance() {
47         if (null == instance) {
48             instance = new NodeDirectory();
49             plugin = plugin.toLowerCase();
50             try {
51                 instance.delegate = (NodeDirectoryInterface) (Class
52                         .forName("ch.nix.ert.communication." + plugin
53                                 + ".NodeDirectory").newInstance());
54                 return instance;
55             } catch (InstantiationException e) {
56                 // TODO Auto-generated catch block
57                 e.printStackTrace();
58             } catch (IllegalAccessException e) {
59                 // TODO Auto-generated catch block
60                 e.printStackTrace();
61             } catch (ClassNotFoundException e) {
62                 // TODO Auto-generated catch block
63                 e.printStackTrace();
64             }
65             return null;
66         }
67         return instance;
68
69     }
70
71     /*
72      * (non-Javadoc)
73      *
74      * @see ch.nix.ert.communication.NodeDirectoryInterface#publish(ch.nix.ert.KeyPartition)
75      */
76     public void publish(DataItem keyPartition) {
77         logger.info("Publish: "+keyPartition);
78         delegate.publish(keyPartition);
79     }
80
81     /*
82      * (non-Javadoc)
83      *
84      * @see ch.nix.ert.communication.NodeDirectoryInterface#unpublish(ch.nix.ert.KeyPartition)
85      */
86     public void unpublish(DataItem keyPartition) {
87         logger.info("Unpublish: "+keyPartition);
88         delegate.unpublish(keyPartition);
89     }
90
91     /*
92      * (non-Javadoc)
93      *
94      * @see ch.nix.ert.communication.NodeDirectoryInterface#findNode(ch.nix.ert.DataItem)
95      */
96
97     public DataItem findNode(DataItem data) {
98         return delegate.findNode(data);
99     }
100
101 //    public Node findNode(InternalKey key) {
102 //        DataItem myData = new DataItem();
103 //        myData.putValue("internalKey", key);
104 //        return this.findNode(myData);
105 //    }
106
107 }
```

---

**File 12:** communication.NodeDirectoryInterface

---

```
 1 /*
 2  * Created on Aug 8, 2004 8:12:11 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 138 $
 8  * Location:    $URL: NodeDirectoryInterface.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication;
17
18 import ch.nix.ert.message.DataItem;
19
20 /**
21  * TODO Edit this text for a better description
22  *
23  * @author michi
24  */
25 public interface NodeDirectoryInterface {
26
27     public void publish(DataItem dataKeyPartition);
28     public void unpublish(DataItem dataKeyPartition);
29     public DataItem findNode(DataItem data);
30
31 }
```

---

**File 13:** communication.jmdns.NodeDirectory

---

```
 1 /*
 2  * Created on Aug 4, 2004 3:29:26 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-11-11 19:33:00 +0100 (Thu, 11 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 171 $
 8  * Location:    $URL: NodeDirectory.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.jmdns;
17
18 import java.io.IOException;
19 import java.util.Hashtable;
20 import java.util.SortedMap;
21
22 import javax.jmdns.JmDNS;
23 import javax.jmdns.ServiceInfo;
24
25 import org.apache.log4j.Logger;
26
27 import ch.nix.ert.communication.NodeDirectoryInterface;
28 import ch.nix.ert.message.DataItem;
29 import ch.nix.ert.message.InternalKey;
30 import ch.nix.ert.message.Node;
31
32 /**
33  * TODO Edit this text for a better description
34  *
35  * This is a special implementation for zeroconf (aka Apple:Rendezvous,
36  * Apple:OpenTalk) and published all KeyPartitions
37  *
38  * @author michi
39  */
40 public class NodeDirectory implements NodeDirectoryInterface {
41
42     // logger:
43     static Logger logger = Logger.getLogger(NodeDirectory.class);
44
45     private JmDNS jmdns;
```

```
46
47     private Hashtable registeredServiceTable = new Hashtable();
48
49     private NodeDirectoryListener nodeDirectoryListener = new NodeDirectoryListener();
50
51     /**
52      *
53      */
54     public NodeDirectory() {
55         try {
56             jmdns = new JmDNS();
57             //now start Listener to all ExtendedRootTree nodes out there!
58             jmdns.addServiceListener("_master._ert._tcp.local.",
59                     nodeDirectoryListener);
60         } catch (IOException e) {
61             // TODO Auto-generated catch block
62             logger.error("Error creating zeroconf service (JmDNS)" + e);
63         }
64     }
65
66     // TODO change: the HashTable should contain a reference to the actual
67     // KeyPartition, and should check if the information in this keyPartition
68     // has changed: if so, it has to unpublish the old and publish the new one
69     public void publish(DataItem data) {
70         //KeyPartition keyPartition = data.getKeyPartition();
71         InternalKey firstInternalKey = data.getFirstInternalKey();
72         InternalKey lastInternalKey = data.getLastInternalKey();
73         int keyPartitionId = data.getKeyPartitionId();
74
75         Hashtable props = new Hashtable();
76         //now check, if for some reason something is alreay registered at the
77         // very same name! (To prevent the creation of [2]...)
78         if (nodeDirectoryListener.getKeyPartitionList().containsKey(
79                 firstInternalKey)) {
80             //there is already something out there with the name 2, I dont
81             // publish this new information!
82             //System.out.println("Already something published with this
83             // name!");
84             logger
85                     .info("There is already something published with this KeyGroup-ID: "
86                             + firstInternalKey
87                             + ", I don't overwrite this informatin");
88             return;
89         }
90         String publishedInfo = firstInternalKey.toString();
91         if (registeredServiceTable.containsKey(new Integer(keyPartitionId))) {
92             //This service is already published, now check, if KeyPartition has
93             // changed!
94             ServiceInfo oldServiceInfo = (ServiceInfo) registeredServiceTable
95                     .get(new Integer(keyPartitionId));
96             String oldPublishedInfo = oldServiceInfo.getName();
97             logger.info("KeyPartition was already published, but has changed"
98                     + " in meantime, I will update my Directory from: "
99                     + oldPublishedInfo + " to: " + publishedInfo);
100            //System.out.println("new: " + publishedInfo);
101            //System.out.println("old: " + oldPublishedInfo);
102            //System.out.println("Same?" +
103            // publishedInfo.equals(oldPublishedInfo));
104            if (publishedInfo.equals(oldPublishedInfo)) {
105                // keyPartition is the same, nothing to do!
106                //System.out.println("already published, nothing todo!");
107                logger
108                        .info("KeyPartition has changed, but is not relevant"
109                                + " for NodeDirectory, so I don't re-publish anything now");
110                return;
111            } else {
112                // keyPartitions have significantly changed, i need to update my
113                // directory
114                //System.out.println("unpublish/publish me: " + publishedInfo);
115                logger
116                        .info("KeyPartition has changed, and changes are relevant"
117                                + " for NodeDirectory, I will now re-publish"
118                                + " the updated information");
119                jmdns.unregisterService(oldServiceInfo);
120                registeredServiceTable.remove(new Integer(keyPartitionId));
121            }
122        }
123        props.put("keygroup", publishedInfo);
124        // TODO currently the port 1099 is hardcoded (RMI port), must be made
125        // configurable.
126        ServiceInfo newServiceInfo = new ServiceInfo(
127                "_master._ert._tcp.local.", publishedInfo
128                        + "._master._ert._tcp.local.", 1099, 0, 0, props);
129
130        try {
```

```
131              //System.out.println("registering new service");
132              logger.info("Register KeyPartition: (" + firstInternalKey + "-"
133                      + lastInternalKey + ") with ID: " + keyPartitionId
134                      + " (JmDNS ServiceInfo: "
135                      + newServiceInfo.getNiceTextString() + ")");
136              jmdns.registerService(newServiceInfo);
137              registeredServiceTable.put(new Integer(keyPartitionId),
138                      newServiceInfo);
139          } catch (IOException e) {
140              // TODO Auto-generated catch block
141              System.out.println("ERROR registering zeroconf service (JmDNS) "
142                      + e);
143              logger
144                      .error("Error registering KeyPartition in zeroconf service (JmDNS)");
145          }
146      }
147
148      public void unpublish(DataItem data) {
149          //KeyPartition keyPartition = data.getKeyPartition();
150          InternalKey firstInternalKey = data.getFirstInternalKey();
151          InternalKey lastInternalKey = data.getLastInternalKey();
152          int keyPartitionId = data.getKeyPartitionId();
153
154          ServiceInfo oldServiceInfo = (ServiceInfo) registeredServiceTable
155                  .get(new Integer(keyPartitionId));
156          //System.out.println("KeyPartition: " + keyPartition);
157          //System.out.println("ServiceHash: " + registeredServiceTable);
158          if (oldServiceInfo != null) {
159              logger.info("Unregister KeyPartition with ID: " + keyPartitionId
160                      + " (JmDNS ServiceInfo: "
161                      + oldServiceInfo.getNiceTextString() + ")");
162              jmdns.unregisterService(oldServiceInfo);
163          }
164          registeredServiceTable.remove(new Integer(keyPartitionId));
165      }
166
167      public DataItem findNode(DataItem data) {
168          if (data.getInternalKey() == null) {
169              return null;
170          }
171          SortedMap headMap = nodeDirectoryListener.getKeyPartitionList()
172                  .headMap(data.getInternalKey().getSubsequentInternalKey());
173          //if this beast is not empty, get the last key of it:
174          Node node;
175          if (!headMap.isEmpty()) {
176              node = (Node) nodeDirectoryListener.getKeyPartitionList().get(
177                      headMap.lastKey());
178              DataItem returnData = new DataItem();
179              returnData.setDirectDestination(node);
180              return returnData;
181          }
182          // should not happen, as always all KeyPartitions should be available
183          // (only in case a machine is down, network problem, this can happen
184          return null;
185      }
186 }
```

---

**File 14:** communication.jmdns.NodeDirectoryListener

```
1 /*
2  * Created on Aug 4, 2004 5:31:31 PM
3  *
4  * Project:      ExtendendRootTree
5  * Last Change: $Date: 2004-11-04 10:15:02 +0100 (Thu, 04 Nov 2004) $
6  * Changed by:   $Author: michi $
7  * Revision:     $Rev: 158 $
8  * Location:     $URL: NodeDirectoryListener.java $
9  *
10 * This source-code is part of the diploma thesis of Michael Kussmaul
11 * written at the University of Zurich, Switzerland
12 * Department of Information Technology (www.ifi.unizh.ch)
13 * For copyright information please contact kussmaul@nix.ch
14 *
15 */
16 package ch.nix.ert.communication.jmdns;
17
18 import java.util.TreeMap;
19
20 import javax.jmdns.JmDNS;
21 import javax.jmdns.ServiceInfo;
22 import javax.jmdns.ServiceListener;
23
```

```
24 import ch.nix.ert.message.InternalKey;
25 import ch.nix.ert.message.Node;
26
27 /**
28  * TODO Edit this text for a better description
29  *
30  * @author michi
31  */
32 public class NodeDirectoryListener implements ServiceListener {
33
34     private TreeMap nodeDirectory = new TreeMap();
35
36     public TreeMap getKeyPartitionList() {
37         return nodeDirectory;
38     }
39
40     /*
41      * (non-Javadoc)
42      *
43      * @see javax.jmdns.ServiceListener#addService(javax.jmdns.JmDNS,
44      *      java.lang.String, java.lang.String)
45      */
46     public void addService(JmDNS jmdns, String type, String name) {
47         String shortname;
48         shortname = name.substring(0, name.length() - (type.length() + 1));
49         //check if name is a number (e.g. not 1 [1])
50         if (shortname.endsWith("]")) {
51             //invalid name, discard it!
52             //System.out.println("Wrong Name in addService: " + shortname);
53             return;
54         }
55         InternalKey key = new InternalKey(shortname);
56         Node value = new Node(jmdns.getServiceInfo(type, name).getAddress(),
57                 jmdns.getServiceInfo(type, name).getPort());
58         nodeDirectory.put(key, value);
59         //System.out.println("ADD: " + shortname);
60     }
61
62     /*
63      * (non-Javadoc)
64      *
65      * @see javax.jmdns.ServiceListener#removeService(javax.jmdns.JmDNS,
66      *      java.lang.String, java.lang.String)
67      */
68     public void removeService(JmDNS jmdns, String type, String name) {
69         // TODO Auto-generated method stub
70         if (name.endsWith("." + type)) {
71             name = name.substring(0, name.length() - (type.length() + 1));
72         }
73         // check if name is a number (e.g. not 1 [1])
74         if (name.endsWith("]")) {
75             //invalid name, discard it!
76             //System.out.println("Wrong Name in removeService: " + name);
77             return;
78         }
79         InternalKey key = new InternalKey(name);
80         nodeDirectory.remove(key);
81         //System.out.println("REMOVE: " + name);
82     }
83
84     /*
85      * (non-Javadoc)
86      *
87      * @see javax.jmdns.ServiceListener#resolveService(javax.jmdns.JmDNS,
88      *      java.lang.String, java.lang.String, javax.jmdns.ServiceInfo)
89      */
90     public void resolveService(JmDNS jmdns, String type, String name,
91             ServiceInfo info) {
92         // TODO Auto-generated method stub
93         System.out.println("RESOLVED: " + info);
94     }
95
96 }
```

**File 15:** communication.netty2.ErtClassCatalog

```
1 /*
2  * Created on Oct 22, 2004 9:52:21 PM
3  *
4  * Project:      ExtendendRootTree
5  * Last Change: $Date: 2004-11-12 12:46:25 +0100 (Fri, 12 Nov 2004) $
6  * Changed by:  $Author: michi $
```

```
 7  * Revision:    $Rev: 173 $
 8  * Location:    $URL: ErtClassCatalog.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.netty2;
17
18 import java.io.ObjectStreamClass;
19 import java.util.Hashtable;
20
21 import org.apache.log4j.Logger;
22
23 import ch.nix.ert.message.DataItem;
24
25 import com.sleepycat.bind.serial.ClassCatalog;
26 import com.sleepycat.je.DatabaseException;
27
28 /**
29  * TODO Edit this text for a better description
30  *
31  * @author michi
32  */
33 public class ErtClassCatalog implements ClassCatalog {
34
35     // logger:
36     static Logger logger = Logger.getLogger(ErtClassCatalog.class);
37
38     String[] classIds = { "ch.nix.ert.DataItem", "java.util.Hashtable",
39             "java.lang.Boolean", "ch.nix.ert.Node",
40             "java.net.InetSocketAddress", "java.net.InetAddress",
41             "java.net.SocketAddress", "ch.nix.ert.InternalKey",
42             "ch.nix.ert.UserKey" };
43
44     /*
45      * (non-Javadoc)
46      *
47      * @see com.sleepycat.bind.serial.ClassCatalog#close()
48      */
49     public void close() throws DatabaseException {
50         // TODO Auto-generated method stub
51
52     }
53
54     /*
55      * (non-Javadoc)
56      *
57      * @see com.sleepycat.bind.serial.ClassCatalog#getClassFormat(byte[])
58      */
59     public ObjectStreamClass getClassFormat(byte[] classID)
60             throws DatabaseException, ClassNotFoundException {
61         // System.out.println("getClassFormat:
62         // "+ObjectStreamClass.lookup(DataItem.class));
63         //System.out.println("Lookedup: " + classID[0]);
64         String name = classIds[(int) classID[0]];
65         return ObjectStreamClass.lookup(Class.forName(name));
66     }
67
68     /*
69      * (non-Javadoc)
70      *
71      * @see com.sleepycat.bind.serial.ClassCatalog#getClassID(java.io.ObjectStreamClass)
72      */
73     public byte[] getClassID(ObjectStreamClass classDesc)
74             throws DatabaseException, ClassNotFoundException {
75         // TODO Auto-generated method stub
76         // System.out.println("getClassID: "+classDesc);
77         String name = classDesc.getName();
78         for (int i = 0; i < classIds.length; i++) {
79             if (((String) classIds[i]).equals(name)) {
80                 // System.out.println("calculated: " + i);
81                 byte[] result = new byte[1];
82                 result[0] = (byte) i;
83                 return result;
84             }
85
86         }
87         logger.error("Undefined Class: " + name
88                 + ", please update ErtClassCatalog with the"
89                 + " new Class you introduced in DataItem");
90         return new byte[0xff];
91     }
```

```
92
93 }
```

---

**File 16:** communication.netty2.ErtMessage

```
 1 /*
 2  * Created on Oct 21, 2004 9:24:15 AM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-28 22:38:52 +0200 (Thu, 28 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:     $Rev: 144 $
 8  * Location:     $URL: ErtMessage.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.netty2;
17
18 import java.io.ByteArrayInputStream;
19 import java.io.ByteArrayOutputStream;
20 import java.io.IOException;
21 import java.io.ObjectOutputStream;
22 import java.nio.ByteBuffer;
23
24 import net.gleamynode.netty2.Message;
25 import net.gleamynode.netty2.MessageParseException;
26
27 import org.apache.log4j.Logger;
28
29 import ch.nix.ert.message.DataItem;
30 import ch.nix.ert.message.DataItemCoder;
31
32 import com.sleepycat.bind.serial.SerialInput;
33 import com.sleepycat.bind.serial.SerialOutput;
34
35 /**
36  * TODO Edit this text for a better description
37  *
38  * @author michi
39  */
40 public class ErtMessage implements Message {
41
42     // logger:
43     static Logger logger = Logger.getLogger(ErtMessage.class);
44
45     public DataItem dataItem = null;
46
47     private DataItemCoder dataItemCoder;
48
49     public ErtMessage() {
50         dataItemCoder = new DataItemCoder();
51     }
52
53     public DataItem getDataItem() {
54         return dataItemCoder.getDataItem();
55     }
56
57     public void setDataItem(DataItem myData) {
58         dataItem = myData;
59         dataItemCoder.setDataItem(myData);
60     }
61
62     /*
63      * (non-Javadoc)
64      *
65      * @see net.gleamynode.netty2.Message#read(java.nio.ByteBuffer)
66      */
67     public boolean read(ByteBuffer buf) throws MessageParseException {
68         return dataItemCoder.decode(buf);
69     }
70
71     /*
72      * (non-Javadoc)
73      *
74      * @see net.gleamynode.netty2.Message#write(java.nio.ByteBuffer)
75      */
76     public boolean write(ByteBuffer buf) {
77         return dataItemCoder.encode(buf);
```

```
78      }
79
80 }
```

---

**File 17:** communication.netty2.ErtMessageRecognizer

```
 1 /*
 2  * Created on Oct 21, 2004 9:23:30 AM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-29 17:46:07 +0200 (Fri, 29 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 146 $
 8  * Location:    $URL: ErtMessageRecognizer.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.netty2;
17
18 import java.nio.ByteBuffer;
19
20 import net.gleamynode.netty2.Message;
21 import net.gleamynode.netty2.MessageParseException;
22 import net.gleamynode.netty2.MessageRecognizer;
23
24 /**
25  * TODO Edit this text for a better description
26  *
27  * @author michi
28  */
29 public class ErtMessageRecognizer implements MessageRecognizer {
30
31     /* (non-Javadoc)
32      * @see net.gleamynode.netty2.MessageRecognizer#recognize(java.nio.ByteBuffer)
33      */
34     public  Message recognize(ByteBuffer buf) throws MessageParseException {
35         // return null if we have not yet received the length of the message
36         if (buf.remaining() < 5) {
37             return null;
38         }
39         // we have the length so we get the message!
40 //        int length = buf.getInt();
41 //        if (length < 1) {
42 //            return null;
43 //        }
44         return new ErtMessage();
45     }
46
47 }
```

---

**File 18:** communication.netty2.NodeCommunication

```
 1 /*
 2  * Created on Oct 21, 2004 8:30:54 AM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-11-09 01:55:31 +0100 (Tue, 09 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 165 $
 8  * Location:    $URL: NodeCommunication.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.netty2;
17
18 import java.io.IOException;
19 import java.net.InetSocketAddress;
20
21 import org.apache.log4j.Logger;
22
23 import net.gleamynode.netty2.IoProcessor;
```

```
24 import net.gleamynode.netty2.LowLatencyEventDispatcher;
25 import net.gleamynode.netty2.MessageRecognizer;
26 import net.gleamynode.netty2.OrderedEventDispatcher;
27 import net.gleamynode.netty2.Session;
28 import net.gleamynode.netty2.SessionServer;
29 import net.gleamynode.netty2.SimpleEventDispatcher;
30 import net.gleamynode.netty2.ThreadPooledEventDispatcher;
31 import ch.nix.ert.base.Dispatcher;
32 import ch.nix.ert.base.ServiceManager;
33 import ch.nix.ert.communication.NodeCommunicationInterface;
34 import ch.nix.ert.communication.NodeDirectory;
35 import ch.nix.ert.communication.rmi.NodeCommunicationRmiInterface;
36 import ch.nix.ert.message.DataItem;
37 import ch.nix.ert.message.Node;
38
39 /**
40  * This is an implementation of Netty 2
41  * (http://gleamynode.net/dev/projects/netty2)
42  *
43  * It is implemented on top of Java NIO and provides good network performance
44  *
45  * @author michi
46  */
47 public class NodeCommunication implements NodeCommunicationInterface {
48
49     // logger:
50     static Logger logger = Logger.getLogger(NodeCommunication.class);
51
52     private static final String HOSTNAME = "localhost";
53
54     private static final int SERVER_PORT = 15123;
55
56     private static final int CONNECT_TIMEOUT = 30;
57
58     private static final int DISPATCHER_THREAD_POOL_SIZE = 1;
59
60     private NodeCommunicationServerListener listener;
61
62     NodeDirectory nodeDirectory = NodeDirectory.getInstance();
63
64     private Dispatcher dispatcher;
65
66     public NodeCommunication() {
67
68         dispatcher = ServiceManager.getInstance().getDispatcher();
69
70         // initialize I/O processor and event dispatcher
71         IoProcessor ioProcessor = new IoProcessor();
72         SimpleEventDispatcher eventDispatcher = new SimpleEventDispatcher();
73
74         try {
75             // start with the default number of I/O worker threads
76             ioProcessor.start();
77         } catch (IOException e) {
78             // TODO Auto-generated catch block
79             e.printStackTrace();
80         }
81
82         // start with a few event dispatcher threads
83         eventDispatcher.setThreadPoolSize(DISPATCHER_THREAD_POOL_SIZE);
84         eventDispatcher.start();
85
86         // prepare message recognizer
87         ErtMessageRecognizer recognizer = new ErtMessageRecognizer();
88
89         // prepare session event listener which will provide communication
90         // workflow.
91         listener = new NodeCommunicationServerListener();
92
93         // prepare session server
94
95         SessionServer server = new SessionServer();
96         server.setIoProcessor(ioProcessor);
97         server.setEventDispatcher(eventDispatcher);
98         server.setMessageRecognizer(recognizer);
99
100         server.addSessionListener(listener);
101         server.setBindAddress(new InetSocketAddress(SERVER_PORT));
102
103         // open the server port, accept connections, and start communication
104
105         // logger.info("Listening on port " + SERVER_PORT);
106         try {
107             server.start();
108         } catch (IOException e1) {
```

```
109                    // TODO Auto-generated catch block
110                    e1.printStackTrace();
111            }
112        }
113
114        /*
115         * (non-Javadoc)
116         *
117         * @see ch.nix.ert.communication.NodeCommunicationInterface#send(ch.nix.ert.DataItem)
118         */
119        public void send(DataItem dataItem) {
120            String ip = "";
121            // now check where the dataItem has to be sent to
122            // first check if we have a "directdelivery" packet, which does not need
123            // to be routed
124            Node node = dataItem.getDirectDestination();
125            if (node == null) {
126                // no indication that we have a direct delivery packet, so we have
127                // to lookup the dataItem
128                DataItem destination = nodeDirectory.findNode(dataItem);
129                if (destination == null) {
130                    //                logger
131                    //                        .warn("Could not send (unknown destination): "
132                    //                                + dataItem);
133
134                    return;
135                }
136                Node destinationNode = destination.getDirectDestination();
137                if (destinationNode.isLocalhost()) {
138                    // means we would send this dataItem to ourselfs, so we return
139                    // it here already!
140                    dataItem.setRouting("local");
141                    dispatcher.addWork(dataItem);
142                    return;
143                }
144                ip = destinationNode.getIp();
145
146            } else {
147                ip = node.getIp();
148                // we don't want to send over this information
149                // TODO later we should write an cleanup method into DataItem which
150                // automatically removes all unused parameters bevore sending it or
151                // storing it...
152                dataItem.removeDirectDestination();
153            }
154
155            //now add originator Address:
156            dataItem.setDirectOriginator(new Node(SERVER_PORT));
157
158            Session session = getSessionForAddress(ip);
159            ErtMessage ertMessage = new ErtMessage();
160            ertMessage.setDataItem(dataItem);
161            //        logger.debug("for IP: " + ip + ", I took session: "
162            //                + session.getSocketAddressString() + " " + session);
163            //        session.setAttachment(this);
164            //        boolean result = session.write(ertMessage);
165            boolean result = false;
166            int round = 0;
167            while (!result && round < 2) {
168                // sending was not successful, try again!
169                session = getSessionForAddress(ip);
170                result = session.write(ertMessage);
171                round++;
172            }
173        }
174
175        private Session getSessionForAddress(String ip) {
176            Session cachedSession = null;
177            cachedSession = SessionManager.getSessionForAddress(ip);
178            //System.out.println(SessionManager.getSize());
179            if (cachedSession == null) {
180                logger
181                        .info("In hashtable are "
182                                + SessionManager.getSize()
183                                + " connections stored, I need to create a new session for ip: "
184                                + ip);
185                IoProcessor clientIoProcessor = new IoProcessor();
186                try {
187                    clientIoProcessor.start();
188                } catch (IOException e) {
189                    // TODO Auto-generated catch block
190                    e.printStackTrace();
191                }
192                SimpleEventDispatcher clientEventDispatcher = new SimpleEventDispatcher();
193                clientEventDispatcher
```

```
194                     .setThreadPoolSize(DISPATCHER_THREAD_POOL_SIZE);
195             clientEventDispatcher.start();
196             ErtMessageRecognizer clientRecognizer = new ErtMessageRecognizer();
197             NodeCommunicationClientListener clientListener = new NodeCommunicationClientListener();
198
199             cachedSession = new Session(clientIoProcessor,
200                     new InetSocketAddress(ip, SERVER_PORT), clientRecognizer,
201                     clientEventDispatcher);
202             // set configuration
203             cachedSession.getConfig().setConnectTimeout(CONNECT_TIMEOUT);
204             // suscribe and start communication
205             cachedSession.addSessionListener(clientListener);
206             SessionManager.addAddressForSession(
207                     ((InetSocketAddress) cachedSession.getSocketAddress())
208                             .getAddress().getHostAddress(), cachedSession);
209             //          logger.debug("Connecting to "
210             //                  + ((InetSocketAddress) cachedSession.getSocketAddress())
211             //                          .getAddress().getHostAddress());
212             cachedSession.start();
213             try {
214                 while (!cachedSession.isConnected()) {
215                     Thread.sleep(500);
216                 }
217             } catch (InterruptedException e1) {
218                 // TODO Auto-generated catch block
219                 e1.printStackTrace();
220                 logger.error("Interrupt: " + e1);
221             }
222         }
223         return cachedSession;
224     }
225
226 }
```

---

**File 19:** communication.netty2.NodeCommunicationClientListener

```
 1 /*
 2  * Created on Oct 21, 2004 7:36:59 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-11-12 12:46:25 +0100 (Fri, 12 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 173 $
 8  * Location:    $URL: NodeCommunicationClientListener.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.netty2;
17
18 import java.net.InetSocketAddress;
19
20 import net.gleamynode.netty2.Message;
21 import net.gleamynode.netty2.Session;
22 import net.gleamynode.netty2.SessionListener;
23
24 import org.apache.log4j.Logger;
25
26 import ch.nix.ert.base.Dispatcher;
27 import ch.nix.ert.base.ServiceManager;
28 import ch.nix.ert.message.DataItem;
29
30 /**
31  * TODO Edit this text for a better description
32  *
33  * @author michi
34  */
35 public class NodeCommunicationClientListener implements SessionListener {
36
37     // logger:
38     static Logger logger = Logger
39             .getLogger(NodeCommunicationClientListener.class);
40
41     private Dispatcher dispatcher;
42
43     public NodeCommunicationClientListener() {
44         dispatcher = ServiceManager.getInstance().getDispatcher();
45     }
46
```

```
47      /*
48       * (non-Javadoc)
49       *
50       * @see net.gleamynode.netty2.SessionListener#
51       * connectionEstablished(net.gleamynode.netty2.Session)
52       */
53      public void connectionEstablished(Session session) {
54          logger.debug("netty2 connection established: " + session);
55
56          // set idle time to 60 seconds
57          session.getConfig().setIdleTime(60);
58          SessionManager.addAddressForSession(((InetSocketAddress) session
59                  .getSocketAddress()).getAddress().getHostAddress(), session);
60      }
61
62      /*
63       * (non-Javadoc)
64       *
65       * @see net.gleamynode.netty2.SessionListener#connectionClosed(net.gleamynode.netty2.Session)
66       */
67      public void connectionClosed(Session session) {
68          logger.debug("netty2 connection closed: " + session);
69          SessionManager.removeAddressForSession(((InetSocketAddress) session
70                  .getSocketAddress()).getAddress().getHostAddress());
71      }
72
73      /*
74       * (non-Javadoc)
75       *
76       * @see net.gleamynode.netty2.SessionListener#messageReceived(net.gleamynode.netty2.Session,
77       *      net.gleamynode.netty2.Message)
78       */
79      public void messageReceived(Session session, Message message) {
80          DataItem dataItem = ((ErtMessage) message).getDataItem();
81          // logger.debug("Recveived <- : " + dataItem);
82          dataItem.setRouting("local");
83          dispatcher.addWork(dataItem);
84      }
85
86      /*
87       * (non-Javadoc)
88       *
89       * @see net.gleamynode.netty2.SessionListener#messageSent(net.gleamynode.netty2.Session,
90       *      net.gleamynode.netty2.Message)
91       */
92      public void messageSent(Session session, Message message) {
93          //      DataItem dataItem = ((ErtMessage) message).getDataItem();
94          //      logger.debug("Sending -> "
95          //              + ((InetSocketAddress) session.getSocketAddress()).getAddress()
96          //                      .getHostAddress() + " : " + dataItem);
97
98      }
99
100     /*
101      * (non-Javadoc)
102      *
103      * @see net.gleamynode.netty2.SessionListener#sessionIdle(net.gleamynode.netty2.Session)
104      */
105     public void sessionIdle(Session session) {
106         //  logger.info("Disconnecting the idle session: " + session);
107
108         SessionManager.removeAddressForSession(((InetSocketAddress) session
109                 .getSocketAddress()).getAddress().getHostAddress());
110         session.close();
111     }
112
113     /*
114      * (non-Javadoc)
115      *
116      * @see net.gleamynode.netty2.SessionListener#exceptionCaught(net.gleamynode.netty2.Session,
117      *      java.lang.Throwable)
118      */
119     public void exceptionCaught(Session session, Throwable cause) {
120         // cause.printStackTrace();
121         logger.error("Unexpected exception: " + cause + "\n" + " in session: "
122                 + session);
123         SessionManager.removeAddressForSession(((InetSocketAddress) session
124                 .getSocketAddress()).getAddress().getHostAddress());
125         session.close();
126
127     }
128
129 }
```

**File 20:** communication.netty2.NodeCommunicationServerListener

```
 1 /*
 2  * Created on Oct 21, 2004 10:20:43 AM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-11-12 12:46:25 +0100 (Fri, 12 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 173 $
 8  * Location:    $URL: NodeCommunicationServerListener.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.netty2;
17
18 import java.net.ConnectException;
19 import java.net.InetSocketAddress;
20
21 import org.apache.log4j.Logger;
22
23 import com.sleepycat.util.RuntimeExceptionWrapper;
24
25 import ch.nix.ert.base.Dispatcher;
26 import ch.nix.ert.base.ServiceManager;
27 import ch.nix.ert.message.DataItem;
28 import net.gleamynode.netty2.Message;
29 import net.gleamynode.netty2.Session;
30 import net.gleamynode.netty2.SessionListener;
31
32 /**
33  * TODO Edit this text for a better description
34  *
35  * @author michi
36  */
37 public class NodeCommunicationServerListener implements SessionListener {
38
39     // logger:
40     static Logger logger = Logger
41             .getLogger(NodeCommunicationServerListener.class);
42
43     private Dispatcher dispatcher;
44
45     public NodeCommunicationServerListener() {
46         dispatcher = ServiceManager.getInstance().getDispatcher();
47     }
48
49     /*
50      * (non-Javadoc)
51      *
52      * @see net.gleamynode.netty2.SessionListener#
53      * connectionEstablished(net.gleamynode.netty2.Session)
54      */
55     public void connectionEstablished(Session session) {
56         logger.debug("netty2 connection established: " + session);
57         // set idle time to 60 seconds
58         session.getConfig().setIdleTime(60);
59         SessionManager.addAddressForSession(((InetSocketAddress) session
60                 .getSocketAddress()).getAddress().getHostAddress(), session);
61     }
62
63     /*
64      * (non-Javadoc)
65      *
66      * @see net.gleamynode.netty2.SessionListener#
67      * connectionClosed(net.gleamynode.netty2.Session)
68      */
69     public void connectionClosed(Session session) {
70         logger.debug("netty2 connection closed: " + session);
71         SessionManager.removeAddressForSession(((InetSocketAddress) session
72                 .getSocketAddress()).getAddress().getHostAddress());
73     }
74
75     /*
76      * (non-Javadoc)
77      *
78      * @see net.gleamynode.netty2.SessionListener#
79      * messageReceived(net.gleamynode.netty2.Session,
80      *      net.gleamynode.netty2.Message)
81      */
82     public void messageReceived(Session session, Message message) {
```

```
 83            DataItem dataItem = ((ErtMessage) message).getDataItem();
 84            // logger.debug("Recveived <- : " + dataItem);
 85            dataItem.setRouting("local");
 86            dispatcher.addWork(dataItem);
 87        }
 88
 89        /*
 90         * (non-Javadoc)
 91         *
 92         * @see net.gleamynode.netty2.SessionListener#
 93         * messageSent(net.gleamynode.netty2.Session,
 94         *      net.gleamynode.netty2.Message)
 95         */
 96        public void messageSent(Session session, Message message) {
 97            //        DataItem dataItem = ((ErtMessage) message).getDataItem();
 98            //        logger.debug("Sending -> "
 99            //              + ((InetSocketAddress) session.getSocketAddress()).getAddress()
100            //                    .getHostAddress() + " : " + dataItem);
101
102        }
103
104        /*
105         * (non-Javadoc)
106         *
107         * @see net.gleamynode.netty2.SessionListener#sessionIdle(net.gleamynode.netty2.Session)
108         */
109        public void sessionIdle(Session session) {
110            logger.info("Disconnecting the idle session: " + session);
111            SessionManager.removeAddressForSession(((InetSocketAddress) session
112                    .getSocketAddress()).getAddress().getHostAddress());
113            session.close();
114        }
115
116        /*
117         * (non-Javadoc)
118         *
119         * @see net.gleamynode.netty2.SessionListener#
120         * exceptionCaught(net.gleamynode.netty2.Session,
121         *      java.lang.Throwable)
122         */
123        public void exceptionCaught(Session session, Throwable cause) {
124            cause.printStackTrace();
125            logger.error("Unexpected exception: " + cause + "\n" + " in session: "
126                    + session);
127            if (cause instanceof RuntimeExceptionWrapper) {
128                Throwable causeDetails = ((RuntimeExceptionWrapper) cause)
129                        .getDetail();
130                logger.error("  Details: " + causeDetails);
131                logger.error("  Message: " + causeDetails.getMessage());
132            }
133
134            if (cause instanceof ConnectException) {
135                // failed to connect to the server, retry after a while.
136                try {
137                    Thread.sleep(5000);
138                } catch (InterruptedException e) {
139                }
140                session.start();
141            } else {
142                SessionManager.removeAddressForSession(((InetSocketAddress) session
143                        .getSocketAddress()).getAddress().getHostAddress());
144                session.close();
145            }
146
147        }
148
149 }
```

---

**File 21:** communication.netty2.SessionManager

```
 1 /*
 2  * Created on Oct 21, 2004 1:39:20 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 138 $
 8  * Location:    $URL: SessionManager.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
```

```
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.netty2;
17
18 import java.util.Hashtable;
19
20 import ch.nix.ert.message.Node;
21
22 import net.gleamynode.netty2.Session;
23
24 /**
25  * TODO Edit this text for a better description
26  *
27  * @author michi
28  */
29 public class SessionManager {
30
31     private static Hashtable sessionHash = new Hashtable();
32
33     public static Session getSessionForAddress(String ip) {
34         return (Session) sessionHash.get(ip);
35     }
36
37     public static void addAddressForSession(String ip, Session session) {
38         sessionHash.put(ip, session);
39     }
40
41     public static void removeAddressForSession(String ip) {
42         sessionHash.remove(ip);
43     }
44
45     public static int getSize() {
46         return sessionHash.size();
47     }
48
49
50 }
```

### File 22: communication.rmi.NodeCommunication

```
 1 /*
 2  * Created on Aug 5, 2004 3:28:51 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 138 $
 8  * Location:    $URL: NodeCommunication.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.rmi;
17
18 import java.net.MalformedURLException;
19 import java.rmi.Naming;
20 import java.rmi.NotBoundException;
21 import java.rmi.RemoteException;
22 import java.rmi.registry.LocateRegistry;
23 import java.rmi.registry.Registry;
24 import java.rmi.server.ExportException;
25
26 import org.apache.log4j.Logger;
27
28 import ch.nix.ert.base.Dispatcher;
29 import ch.nix.ert.base.ServiceManager;
30 import ch.nix.ert.communication.NodeCommunicationInterface;
31 import ch.nix.ert.communication.NodeDirectory;
32 import ch.nix.ert.message.DataItem;
33 import ch.nix.ert.message.Node;
34
35 /**
36  * TODO Edit this text for a better description
37  *
38  * @author michi
39  */
40 // To do the RMI stub building, do:
41 // cd /Users/michi/Documents/eclipse/workspace/ExtendendRootTree/bin
```

```
42 // rmic ch.nix.ert.communication.rmi.NodeCommunicationRmi
43 public class NodeCommunication implements NodeCommunicationInterface {
44
45     static Logger logger = Logger.getLogger(NodeCommunication.class);
46
47     NodeCommunicationRmi communicationObject;
48
49     NodeDirectory nodeDirectory = NodeDirectory.getInstance();
50
51     private Dispatcher dispatcher;
52
53     /**
54      *
55      */
56     public NodeCommunication() {
57         Registry reg;
58         try {
59             try {
60                 reg = LocateRegistry.createRegistry(1099); // registry
61             } catch (ExportException e) {
62                 reg = LocateRegistry.getRegistry(1099);
63             }
64             dispatcher = ServiceManager.getInstance().getDispatcher();
65             communicationObject = new NodeCommunicationRmi(this);
66             communicationObject.startServer();
67         } catch (RemoteException e) {
68             // TODO Auto-generated catch block
69             System.out.println("Generation of RMI failed" + e);
70         }
71     }
72
73     public void send(DataItem dataItem) {
74         String ip = "";
75         // now check where the dataItem has to be sent to
76         // first check if we have a "directdelivery" packet, which does not need
77         // to be routed
78         Node node = dataItem.getDirectDestination();
79         if (node == null) {
80             // no indication that we have a direct delivery packet, so we have
81             // to lookup the dataItem
82             DataItem destination = nodeDirectory.findNode(dataItem);
83             if (destination == null) {
84                 logger.warn("Could not send (unknown destination): "+dataItem);
85                 return;
86             }
87             Node destinationNode = destination.getDirectDestination();
88             if (destinationNode.isLocalhost()) {
89                 // means we would send this dataItem to ourselfs, so we return it here already!
90                 dataItem.setRouting("local");
91                 dispatcher.addWork(dataItem);
92                 return;
93             }
94             ip = destinationNode.getIp();
95
96         } else {
97             ip = node.getIp();
98             // we don't want to send over this information
99             // TODO later we should write an cleanup method into DataItem which
100            // automatically removes all unused parameters bevore sending it or
101            // storing it...
102            dataItem.removeDirectDestination();
103        }
104
105        NodeCommunicationRmiInterface nodeCommunication;
106
107        //now add originator Address:
108        dataItem.setDirectOriginator(new Node(1099));
109
110        try {
111            nodeCommunication = (NodeCommunicationRmiInterface) Naming
112                    .lookup("rmi://" + ip + "/NodeCommunicationRmi");
113            // executes the remote "receive" method to acutally send this dataItem
114            //System.out.println("Data sendingRMI ->: "+ dataItem);
115            logger.debug("Sending -> "+ip+" : "+dataItem);
116            nodeCommunication.receive(dataItem);
117        } catch (MalformedURLException e) {
118            logger.warn("Could not send: "+dataItem+"\n"+e);
119        } catch (RemoteException e) {
120            logger.warn("Could not send: "+dataItem+"\n"+e);
121        } catch (NotBoundException e) {
122            logger.warn("Could not send: "+dataItem+"\n"+e);
123        }
124    }
125 }
```

**File 23:** communication.rmi.NodeCommunicationRmi

```
1  /*
2   * Created on Aug 5, 2004 4:20:49 PM
3   *
4   * Project:       ExtendendRootTree
5   * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
6   * Changed by:  $Author: michi $
7   * Revision:     $Rev: 138 $
8   * Location:     $URL: NodeCommunicationRmi.java $
9   *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.rmi;
17
18 import java.net.MalformedURLException;
19 import java.rmi.Naming;
20 import java.rmi.RemoteException;
21 import java.rmi.server.UnicastRemoteObject;
22
23 import org.apache.log4j.Logger;
24
25 import ch.nix.ert.base.Dispatcher;
26 import ch.nix.ert.base.ServiceManager;
27 import ch.nix.ert.message.DataItem;
28
29 /**
30  * TODO Edit this text for a better description
31  *
32  * @author michi
33  */
34 public class NodeCommunicationRmi extends UnicastRemoteObject implements
35         NodeCommunicationRmiInterface {
36
37     static Logger logger = Logger.getLogger(NodeCommunicationRmi.class);
38
39     NodeCommunication nodeCommunication;
40
41     private Dispatcher dispatcher;
42
43     /**
44      * @throws RemoteException
45      */
46     public NodeCommunicationRmi(NodeCommunication myNodeCommunication)
47             throws RemoteException {
48         nodeCommunication = myNodeCommunication;
49         //write the responses directly to the Dispatcher main queue:
50         dispatcher = ServiceManager.getInstance().getDispatcher();
51     }
52
53     /*
54      * (non-Javadoc)
55      *
56      * @see ch.nix.ert.NodeCommunicationRmiInterface#transmit(ch.nix.ert.DataItem)
57      */
58     public void receive(DataItem dataItem) throws RemoteException {
59         // This will be executed, when remote site was sending us an object
60         //System.out.println("Comm Received: "+dataItem);
61         logger.debug("Recveived <- : "+dataItem);
62         dataItem.setRouting("local");
63         dispatcher.addWork(dataItem);
64     }
65
66     public void startServer() {
67
68         try {
69             Naming.rebind("NodeCommunicationRmi", NodeCommunicationRmi.this);
70             //System.out.println ("NodeCommunication Server is ready.");
71         } catch (MalformedURLException e) {
72             // TODO Auto-generated catch block
73             System.out.println("Error in NodeCommunicationRmi, Naming.bind: "
74                     + e);
75         } catch (RemoteException e) {
76             // TODO Auto-generated catch block
77             System.out.println("Error in NodeCommunicationRmi, Naming.bind: "
78                     + e);
79         }
80     }
81 }
```

**File 24:** communication.rmi.NodeCommunicationRmiInterface

```
 1 /*
 2  * Created on Aug 5, 2004 4:21:16 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 138 $
 8  * Location:    $URL: NodeCommunicationRmiInterface.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.communication.rmi;
17
18 import java.rmi.Remote;
19 import java.rmi.RemoteException;
20
21 import ch.nix.ert.message.DataItem;
22
23 /**
24  * TODO Edit this text for a better description
25  *
26  * @author michi
27  */
28 public interface NodeCommunicationRmiInterface extends Remote {
29
30     public void receive(DataItem dataItem) throws RemoteException;
31
32 }
```

**File 25:** datastore.DataStore

```
 1 /*
 2  * Created on Aug 10, 2004 12:28:46 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-11-11 18:53:50 +0100 (Thu, 11 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 168 $
 8  * Location:    $URL: DataStore.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.datastore;
17
18 import ch.nix.ert.message.DataItem;
19
20 /**
21  * TODO Edit this text for a better description
22  *
23  * @author michi
24  */
25 public class DataStore implements DataStoreInterface {
26
27     static private String plugin = "berkeleyje";
28
29     private DataStoreInterface delegate = null;
30
31     public DataStore(String instanceId) {
32         try {
33             delegate = (DataStoreInterface) (Class
34                     .forName("ch.nix.ert.datastore." + plugin + ".DataStore")
35                     .newInstance());
36             delegate.setDatabaseId(instanceId);
37         } catch (InstantiationException e) {
38             // TODO Auto-generated catch block
39             e.printStackTrace();
40         } catch (IllegalAccessException e) {
41             // TODO Auto-generated catch block
42             e.printStackTrace();
43         } catch (ClassNotFoundException e) {
44             // TODO Auto-generated catch block
```

```
 45                    e.printStackTrace();
 46            }
 47        }
 48
 49        /*
 50         * (non-Javadoc)
 51         *
 52         * @see ch.nix.ert.datastore.DataStoreInterface#getDataItem(ch.nix.ert.DataItem)
 53         */
 54        public DataItem getDataItem(DataItem data) {
 55            // TODO Auto-generated method stub
 56            return delegate.getDataItem(data);
 57        }
 58
 59        /*
 60         * (non-Javadoc)
 61         *
 62         * @see ch.nix.ert.datastore.DataStoreInterface#putDataItem(ch.nix.ert.DataItem)
 63         */
 64        public boolean putDataItem(DataItem data) {
 65            // TODO Auto-generated method stub
 66            return delegate.putDataItem(data);
 67        }
 68
 69        /*
 70         * (non-Javadoc)
 71         *
 72         * @see ch.nix.ert.datastore.DataStoreInterface#deleteDataItem(ch.nix.ert.DataItem)
 73         */
 74        public boolean deleteDataItem(DataItem data) {
 75            // TODO Auto-generated method stub
 76            return delegate.deleteDataItem(data);
 77        }
 78
 79        public int getDbSize() {
 80            return delegate.getDbSize();
 81        }
 82
 83        /*
 84         * (non-Javadoc)
 85         *
 86         * @see ch.nix.ert.datastore.DataStoreInterface#getDatabaseId()
 87         */
 88        public String getDatabaseId() {
 89            // TODO Auto-generated method stub
 90            return delegate.getDatabaseId();
 91        }
 92
 93        public void setDatabaseId(String id) {
 94            //do nothing as setting the Database id should not be allowed!
 95        }
 96
 97        /*
 98         * (non-Javadoc)
 99         *
100         * @see ch.nix.ert.datastore.DataStoreInterface#getNextDataItem(ch.nix.ert.DataItem,
101         *      java.lang.String, int)
102         */
103        public DataItem getNextDataItem(DataItem data) {
104            // TODO Auto-generated method stub
105            return delegate.getNextDataItem(data);
106        }
107
108        /*
109         * (non-Javadoc)
110         *
111         * @see ch.nix.ert.datastore.DataStoreInterface#getPreviousDataItem(ch.nix.ert.DataItem,
112         *      java.lang.String, int)
113         */
114        public DataItem getPreviousDataItem(DataItem data) {
115            // TODO Auto-generated method stub
116            return delegate.getPreviousDataItem(data);
117        }
118
119 }
```

**File 26:** datastore.DataStoreInterface

```
1 /*
2  * Created on Aug 10, 2004 12:31:08 PM
3  *
4  * Project:     ExtendendRootTree
```

```
 5  * Last Change: $Date: 2004-10-26 15:43:26 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 139 $
 8  * Location:    $URL: DataStoreInterface.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16  package ch.nix.ert.datastore;
17
18  import ch.nix.ert.message.DataItem;
19
20  /**
21   * TODO Edit this text for a better description
22   *
23   * @author michi
24   */
25  public interface DataStoreInterface {
26
27      public DataItem getNextDataItem(DataItem data);
28
29      public DataItem getPreviousDataItem(DataItem data);
30
31      public DataItem getDataItem(DataItem data);
32
33      public boolean putDataItem(DataItem data);
34
35      public boolean deleteDataItem(DataItem data);
36
37      // public void deleteDataItemRange(DataItem data);
38
39      public int getDbSize();
40
41      /**
42       * If database is still closed, <code>databaseId</code> is null
43       * <p>
44       * If database is opened and running, <code>databaseId</code> contains an
45       * unique Id (namespace like "ch.nix.accountDB")
46       *
47       * @return Returns the databaseId.
48       */
49      public String getDatabaseId();
50
51      public void setDatabaseId(String databaseId);
52
53  }
```

**File 27:** datastore.DataStoreManager

```
 1  /*
 2   * Created on Jul 28, 2004 4:59:10 PM
 3   *
 4   * Project:      ExtendendRootTree
 5   * Last Change: $Date: 2004-10-26 15:43:26 +0200 (Tue, 26 Oct 2004) $
 6   * Changed by:  $Author: michi $
 7   * Revision:    $Rev: 139 $
 8   * Location:    $URL: DataStoreManager.java $
 9   *
10   * This source-code is part of the diploma thesis of Michael Kussmaul
11   * written at the University of Zurich, Switzerland
12   * Department of Information Technology (www.ifi.unizh.ch)
13   * For copyright information please contact kussmaul@nix.ch
14   *
15   */
16
17  package ch.nix.ert.datastore;
18
19  import java.util.Hashtable;
20
21  import ch.nix.ert.datastore.DataStore;
22  import ch.nix.ert.message.DataItem;
23
24  /**
25   *
26   * TODO Edit this text for a better description Manages all DataStore objects.
27   * Each DataStore object is responsible for his own namespace <br>
28   *
29   * Follows the Singelton-Design Pattern after Rodney Waldhoff.
30   * http://radio.weblogs.com/0122027/stories/2003/10/20/implementingTheSingletonPatternInJava.html
```

```
31  *
32  * @author michi
33  */
34 public class DataStoreManager {
35
36     private Hashtable dataStores = new Hashtable();
37
38     static private DataStoreManager instance = null;
39
40     private DataStoreManager() {
41     }
42
43     /**
44      * @return The unique instance of this class.
45      */
46     static public synchronized DataStoreManager getInstance() {
47         if (null == instance) {
48             instance = new DataStoreManager();
49         }
50         return instance;
51     }
52
53     public boolean put(DataItem data) {
54         //check if namespace already exists
55         String myNamespace = (String) data.getNamespace();
56         if (dataStores.get(myNamespace) == null) {
57             //namespace does no exist, I need to create a new DataStore first
58             dataStores.put(myNamespace, new DataStore(myNamespace));
59         }
60         DataStore currentDataStore = (DataStore) dataStores.get(myNamespace);
61         return currentDataStore.putDataItem(data);
62     }
63
64     public DataItem get(DataItem data) {
65         DataStore currentDataStore = (DataStore) dataStores.get(data
66                 .getNamespace());
67         if (currentDataStore == null) {
68             return null;
69         }
70         return currentDataStore.getDataItem(data);
71     }
72
73     public boolean delete(DataItem data) {
74         DataStore currentDataStore = (DataStore) dataStores.get(data
75                 .getNamespace());
76         if (currentDataStore == null) {
77             return false;
78         }
79         return currentDataStore.deleteDataItem(data);
80     }
81
82     public DataItem getNext(DataItem data) {
83         DataStore currentDataStore = (DataStore) dataStores.get(data
84                 .getNamespace());
85         if (currentDataStore == null) {
86             return null;
87         }
88         return currentDataStore.getNextDataItem(data);
89     }
90
91     public DataItem getPrevious(DataItem data) {
92         DataStore currentDataStore = (DataStore) dataStores.get(data
93                 .getNamespace());
94         if (currentDataStore == null) {
95             return null;
96         }
97         return currentDataStore.getPreviousDataItem(data);
98     }
99
100    public int getDbSize(DataItem data) {
101        DataStore currentDataStore = (DataStore) dataStores.get(data
102                .getNamespace());
103        if (currentDataStore == null) {
104            return 0;
105        }
106        return currentDataStore.getDbSize();
107    }
108
109 }
```

**File 28:** datastore.DataStoreService

```
1 /*
```

```
 2  * Created on Aug 13, 2004 3:02:42 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-28 22:38:52 +0200 (Thu, 28 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 144 $
 8  * Location:    $URL: DataStoreService.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.datastore;
17
18 import org.apache.log4j.Logger;
19
20 import ch.nix.ert.base.ServiceBase;
21 import ch.nix.ert.message.DataItem;
22 import ch.nix.ert.message.Node;
23
24 /**
25  * TODO Edit this text for a better description
26  *
27  * @author michi
28  */
29 public class DataStoreService extends ServiceBase {
30
31     // logger:
32     static Logger logger = Logger.getLogger(DataStoreService.class);
33
34     static private DataStoreService instance = null;
35
36     private DataStoreManager dataStoreManager;
37
38     // private Queue serviceQueue = new Queue("DataStoreQueue");
39
40     // private Dispatcher dispatcher;
41
42     private DataStoreService() {
43     }
44
45     static public synchronized DataStoreService getInstance() {
46         if (null == instance) {
47             instance = new DataStoreService();
48             // instance.dispatcher =
49             // ServiceManager.getInstance().getDispatcher();
50             instance.dataStoreManager = DataStoreManager.getInstance();
51             logger.info("DataStoreService started");
52         }
53         return instance;
54     }
55
56     //    private Dispatcher getDispatcher() {
57     //        return dispatcher;
58     //    }
59
60     /*
61      * (non-Javadoc)
62      *
63      * @see ch.nix.ert.ServiceInterface#registerAtDispatcher()
64      */
65     public void registerAtDispatcher() {
66         registerService("putdb");
67         registerService("putdbnoack");
68         registerService("getdb");
69         registerService("getnextdb");
70         registerService("getprevdb");
71         registerService("deletedb");
72         registerService("getsizedb");
73     }
74
75     /*
76      * (non-Javadoc)
77      *
78      * @see java.lang.Runnable#run()
79      */
80     public void run() {
81         //System.out.println("DataStoreService started");
82         Thread.currentThread().setName("ERT.DataStoreService");
83         while (true) {
84             DataItem dataItem = null;
85
86             dataItem = getWork();
```

```
 87
 88                    // System.out.println("got request: " + dataItem);
 89                String operation = dataItem.getOperation();
 90                if (operation.equals("getdb")) {
 91                  //  logger.debug("getDB: " + dataItem);
 92                    Node originator = dataItem.getDirectOriginator();
 93                    DataItem response = dataStoreManager.get(dataItem);
 94                    if (response == null) {
 95                        response = new DataItem();
 96                    }
 97                    response.setOperation("getdb");
 98                    response.setDirectDestination(originator);
 99                    transactionEnd(response);
100                  //  logger.debug("getDBResponse: " + response);
101                } else if (operation.equals("getnextdb")) {
102                  //  logger.debug("getNextDB: " + dataItem);
103                    Node originator = dataItem.getDirectOriginator();
104                    DataItem response = dataStoreManager.getNext(dataItem);
105                    if (response == null) {
106                        response = new DataItem();
107                    }
108                    response.setOperation("getnextdb");
109                    response.setDirectDestination(originator);
110                    transactionEnd(response);
111                  //  logger.debug("getNextDBResponse: " + response);
112                } else if (operation.equals("getprevdb")) {
113                  //  logger.debug("getPrevDB: " + dataItem);
114                    Node originator = dataItem.getDirectOriginator();
115                    DataItem response = dataStoreManager.getPrevious(dataItem);
116                    if (response == null) {
117                        response = new DataItem();
118                    }
119                    response.setOperation("getprevdb");
120                    response.setDirectDestination(originator);
121                    transactionEnd(response);
122                  //  logger.debug("getPrevDBResponse: " + response);
123                } else if (operation.equals("putdb")) {
124                    // logger.debug("putDB: " + dataItem);
125                    boolean outcome = dataStoreManager.put(dataItem);
126                    DataItem response = new DataItem();
127                    response.setResult(outcome);
128                    transactionEnd(response);
129                    // logger.debug("putDBResponse"+outcome+": " + response);
130                } else if (operation.equals("putdbnoack")) {
131                  //  logger.debug("putDB: " + dataItem);
132                    dataStoreManager.put(dataItem);
133                }else if (operation.equals("deletedb")) {
134                  //  logger.debug("deleteDB: " + dataItem);
135                    boolean outcome = dataStoreManager.delete(dataItem);
136                    dataItem.setResult(outcome);
137                    transactionEnd(dataItem);
138                  //  logger.debug("deleteDBResponse: " + dataItem);
139                } else if (operation.equals("getsizedb")) {
140                  //  logger.debug("getSizeDB: " + dataItem);
141                    int size = dataStoreManager.getDbSize(dataItem);
142                    DataItem response = new DataItem();
143                    response.setDbSize(size);
144                    response.setOperation("getsizedb");
145                    transactionEnd(response);
146                  //  logger.debug("getSizeDBResponse: " + response);
147                } else {
148                    logger
149                            .warn("DataItem dropped, because no processing was found for: "
150                                    + dataItem);
151                }
152            }
153        }
154
155 }
```

**File 29:** datastore.berkeleyje.DataStore

```
 1 /*
 2  * Created on Jul 27, 2004 11:49:54 AM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-11-12 12:46:25 +0100 (Fri, 12 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 173 $
 8  * Location:    $URL: DataStore.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
```

```
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16
17 package ch.nix.ert.datastore.berkeleyje;
18
19 import ch.nix.ert.datastore.DataStoreInterface;
20 import ch.nix.ert.message.DataItem;
21 import ch.nix.ert.message.UserKey;
22
23 import com.sleepycat.bind.EntryBinding;
24 import com.sleepycat.bind.serial.SerialBinding;
25 import com.sleepycat.bind.serial.StoredClassCatalog;
26 import com.sleepycat.bind.tuple.TupleBinding;
27 import com.sleepycat.je.Cursor;
28 import com.sleepycat.je.DatabaseException;
29 import com.sleepycat.je.Database;
30 import com.sleepycat.je.DatabaseEntry;
31 import com.sleepycat.je.DatabaseStats;
32 import com.sleepycat.je.LockMode;
33 import com.sleepycat.je.OperationStatus;
34 import com.sleepycat.je.DatabaseConfig;
35 import com.sleepycat.je.Environment;
36 import com.sleepycat.je.EnvironmentConfig;
37 import com.sleepycat.je.StatsConfig;
38
39 import java.io.File;
40
41 import org.apache.log4j.Logger;
42
43 /**
44  * @author michi
45  *
46  * TODO Edit this text for a better description Responsible for a namespace.
47  * TODO Extend this class to allow range queries and browsing!
48  * <p>
49  * During first performance testing, this class should allow to store more than
50  * 4000 objects per second (variance around 4000 - 8000) Each namespace holds
51  * it's own database
52  */
53 public class DataStore implements DataStoreInterface {
54
55     // logger:
56     static Logger logger = Logger.getLogger(DataStore.class);
57
58     private String databaseId;
59
60     private int syncCounter = 0;
61
62     //TODO make configurable in a later version
63     private String databasePath = "/Users/tutor/ert/db";
64
65     private int requestsSync = 50;
66
67     //private Hashtable cursorTable = new Hashtable();
68
69     Environment myDbEnvironment;
70
71     Database myDatabase;
72
73     //EntryBinding myDataBinding;
74
75     EntryBinding myKeyBinding;
76
77     /**
78      *
79      */
80     public DataStore() {
81         super();
82     }
83
84     public void openDatabase() {
85         try {
86             // Open the environment. Create it if it does not alreadyexist.
87             EnvironmentConfig envConfig = new EnvironmentConfig();
88             envConfig.setAllowCreate(true);
89             myDbEnvironment = new Environment(new File(databasePath), envConfig);
90
91             // Open the database. Create it if it does not already exist.
92             DatabaseConfig dbConfig = new DatabaseConfig();
93             // set the Comparator Method:
94             //        dbConfig.setBtreeComparator(Long.class);
95             //        dbConfig.setDuplicateComparator(Long.class);
```

```
96                  dbConfig.setAllowCreate(true);
97                  dbConfig.setSortedDuplicates(false);
98                  myDatabase = myDbEnvironment.openDatabase(null, this
99                          .getDatabaseId(), dbConfig);
100
101                  // Open the database that you use to store your class information.
102                  //myClassDb = myDbEnvironment.openDatabase(null, "classDb",
103                  // dbConfig);
104
105                  // Instantiate the class catalog
106                  //StoredClassCatalog classCatalog = new
107                  // StoredClassCatalog(myClassDb);
108
109                  // Create the binding
110                  //myDataBinding = new SerialBinding(classCatalog, DataItem.class);
111                  myKeyBinding = TupleBinding.getPrimitiveBinding(Long.class);
112                  logger.info("Berkeley Database opened");
113          } catch (DatabaseException dbe) {
114              // Exception handling goes here
115              logger.error("Error opening DB: " + dbe);
116          }
117      }
118
119      public void closeDatabase() {
120          try {
121              //close all open cursors:
122              //              for (Enumeration e = cursorTable.elements();
123              // e.hasMoreElements();) {
124              //                      ((Cursor) e.nextElement()).close();
125              //              }
126
127              if (myDatabase != null) {
128                  myDatabase.close();
129              }
130              //              if (myClassDb != null) {
131              //                  myClassDb.close();
132              //              }
133              if (myDbEnvironment != null) {
134                  myDbEnvironment.sync();
135                  myDbEnvironment.close();
136              }
137              logger.info("Berkeley Database closed");
138          } catch (DatabaseException dbe) {
139              // Exception handling goes here
140              logger.error("Error closing DB: " + dbe);
141          }
142
143      }
144
145      /*
146       * (non-Javadoc)
147       *
148       * @see java.lang.Object#finalize()
149       */
150      protected void finalize() throws Throwable {
151          try {
152              closeDatabase(); // close open database
153          } finally {
154              super.finalize();
155          }
156      }
157
158      public DataItem getNextDataItem(DataItem data) {
159
160          int amount = data.getCursorSize();
161          UserKey key = data.getUserKey();
162
163          //check, if connection to database is already open
164          if (myDatabase == null) {
165              this.openDatabase();
166          }
167          try {
168              // first create a cursor
169              Cursor cursor = myDatabase.openCursor(null, null);
170
171              // Create a pair of DatabaseEntry objects. theKey
172              // is used to perform the search. theData is used
173              // to store the data returned by the get() operation.
174              DatabaseEntry theKey = new DatabaseEntry();
175              DatabaseEntry theData = new DatabaseEntry();
176
177              myKeyBinding.objectToEntry(new Long(key.getId()), theKey);
178
179              if (cursor.getSearchKeyRange(theKey, theData, LockMode.DEFAULT)
180                      == OperationStatus.SUCCESS) {
```

```
181                     // Recreate the data & key
182                     //          DataItem retrievedData = (DataItem) myDataBinding
183                     //                        .entryToObject(theData);
184                 DataItem retrievedData = new DataItem();
185                 retrievedData.setData(new String(theData.getData()));
186                 UserKey retrievedKey = (UserKey) new UserKey(
187                         (Long) myKeyBinding.entryToObject(theKey));
188
189                 //add the namespace attribute to DataItem
190                 retrievedData.setNamespace(this.getDatabaseId());
191                 retrievedData.setUserKey(retrievedKey);
192
193                 // now check, that we retrieve the next key and not the value
194                 // for the current key
195                 if (key.equals(retrievedKey)) {
196                     //means by coincidence the searching key was exactly a
197                     // matching key in db, so we have to retrieve the next
198                     // key...
199                     OperationStatus retVal = cursor.getNext(theKey, theData,
200                             LockMode.DEFAULT);
201                     if (retVal == OperationStatus.SUCCESS) {
202                         //                       retrievedData = (DataItem) myDataBinding
203                         //                        .entryToObject(theData);
204                         retrievedData.setData(new String(theData.getData()));
205                         retrievedKey = (UserKey) new UserKey(
206                             (Long) myKeyBinding.entryToObject(theKey));
207                         // add the namespace attribute to DataItem
208                         retrievedData.setNamespace(this.getDatabaseId());
209                         retrievedData.setUserKey(retrievedKey);
210                     } else {
211                         return null;
212                     }
213                 }
214
215                 // if user requested more than one dataitem, we have to use
216                 // the cursor functionality and put all remaining data
217                 // into "dataitems" field
218                 for (int i = 1; i < amount; i++) {
219                     DataItem moreRetrievedData;
220                     UserKey moreRetrievedKey;
221                     OperationStatus retVal = cursor.getNext(theKey, theData,
222                             LockMode.DEFAULT);
223                     if (retVal == OperationStatus.SUCCESS) {
224                         //                       moreRetrievedData = (DataItem) myDataBinding
225                         //                                .entryToObject(theData);
226                         moreRetrievedData = new DataItem();
227                         moreRetrievedData
228                                 .setData(new String(theData.getData()));
229                         moreRetrievedKey = (UserKey) new UserKey(
230                             (Long) myKeyBinding.entryToObject(theKey));
231                         // add the namespace attribute to DataItem
232                         moreRetrievedData.setNamespace(this.getDatabaseId());
233                         moreRetrievedData.setUserKey(moreRetrievedKey);
234                         // add the additional dataItem into retreivedData
235                         // from before:
236                         retrievedData.addDataItems(moreRetrievedData);
237                     } else {
238                         // do nothing
239                     }
240
241                 }
242                 cursor.close();
243                 return retrievedData;
244             } else {
245                 //System.out.println("ERROR 1");
246                 return null;
247             }
248         } catch (Exception e) {
249             // TODO Auto-generated catch block
250             // Exception handling goes here
251             logger.error("Get: Error hapened: " + e);
252             return null;
253         }
254     }
255
256     public DataItem getPreviousDataItem(DataItem data) {
257
258         int amount = data.getCursorSize();
259         UserKey key = data.getUserKey();
260
261         //check, if connection to database is already open
262         if (myDatabase == null) {
263             this.openDatabase();
264         }
265         try {
```

```
266             // first create a cursor
267             Cursor cursor = myDatabase.openCursor(null, null);
268
269             // Create a pair of DatabaseEntry objects. theKey
270             // is used to perform the search. theData is used
271             // to store the data returned by the get() operation.
272             DatabaseEntry theKey = new DatabaseEntry();
273             DatabaseEntry theData = new DatabaseEntry();
274
275             myKeyBinding.objectToEntry(new Long(key.getId()), theKey);
276             OperationStatus status = OperationStatus.SUCCESS;
277             if (cursor.getSearchKeyRange(theKey, theData, LockMode.DEFAULT)
278                     != OperationStatus.SUCCESS) {
279                 // means we are at the end of the database, so now get the last
280                 // entry of the whole db
281                 status = cursor.getLast(theKey, theData, LockMode.DEFAULT);
282             } else {
283                 // now we are one item to high, so we have the get the previous
284                 // one...
285                 status = cursor.getPrev(theKey, theData, LockMode.DEFAULT);
286             }
287             if (status == OperationStatus.SUCCESS) {
288                 // Recreate the data & key
289                 //                  DataItem retrievedData = (DataItem) myDataBinding
290                 //                      .entryToObject(theData);
291                 DataItem retrievedData = new DataItem();
292                 retrievedData.setData(new String(theData.getData()));
293                 UserKey retrievedKey = (UserKey) new UserKey(
294                         (Long) myKeyBinding.entryToObject(theKey));
295
296                 //add the namespace attribute to DataItem
297                 retrievedData.setNamespace(this.getDatabaseId());
298                 retrievedData.setUserKey(retrievedKey);
299
300                 // now check, that we retrieve the next key and not the value
301                 // for the current key
302                 if (key.equals(retrievedKey)) {
303                     //means by coincidence the searching key was exactly a
304                     // matching key in db, so we have to retrieve the next
305                     // key...
306                     OperationStatus retVal = cursor.getPrev(theKey, theData,
307                             LockMode.DEFAULT);
308                     if (retVal == OperationStatus.SUCCESS) {
309                         //                      retrievedData = (DataItem) myDataBinding
310                         //                              .entryToObject(theData);
311                         retrievedData.setData(new String(theData.getData()));
312                         retrievedKey = (UserKey) new UserKey(
313                                 (Long) myKeyBinding.entryToObject(theKey));
314                         // add the namespace attribute to DataItem
315                         retrievedData.setNamespace(this.getDatabaseId());
316                         retrievedData.setUserKey(retrievedKey);
317                     } else {
318                         return null;
319                     }
320                 }
321
322                 // if user requested more than one dataitem, we have to use
323                 // the cursor functionality and put all remaining data
324                 // into "dataitems" field
325                 for (int i = 1; i < amount; i++) {
326                     DataItem moreRetrievedData;
327                     UserKey moreRetrievedKey;
328                     OperationStatus retVal = cursor.getPrev(theKey, theData,
329                             LockMode.DEFAULT);
330                     if (retVal == OperationStatus.SUCCESS) {
331                         moreRetrievedData = new DataItem();
332                         //                      moreRetrievedData = (DataItem) myDataBinding
333                         //                              .entryToObject(theData);
334                         moreRetrievedData
335                                 .setData(new String(theData.getData()));
336                         moreRetrievedKey = (UserKey) new UserKey(
337                                 (Long) myKeyBinding.entryToObject(theKey));
338                         // add the namespace attribute to DataItem
339                         moreRetrievedData.setNamespace(this.getDatabaseId());
340                         moreRetrievedData.setUserKey(moreRetrievedKey);
341                         // add the additional dataItem into retrievedData
342                         // from before:
343                         retrievedData.addDataItems(moreRetrievedData);
344                     } else {
345                         // do nothing
346                     }
347
348                 }
349                 cursor.close();
350                 return retrievedData;
```

```
351              } else {
352                  // System.out.println("ERROR 1");
353                  return null;
354              }
355          } catch (Exception e) {
356              // TODO Auto-generated catch block
357              // Exception handling goes here
358              logger.error("Get: Error hapened: " + e);
359              return null;
360          }
361      }
362
363      public DataItem getDataItem(DataItem data) {
364          //extract key from data
365          long key = data.getUserKey().getId();
366          //check, if connection to database is already open
367          if (myDatabase == null) {
368              this.openDatabase();
369          }
370          try {
371              // Create a pair of DatabaseEntry objects. theKey
372              // is used to perform the search. theData is used
373              // to store the data returned by the get() operation.
374              DatabaseEntry theKey = new DatabaseEntry();
375              DatabaseEntry theData = new DatabaseEntry();
376              myKeyBinding.objectToEntry(new Long(key), theKey);
377
378              // Perform the get.
379              if (myDatabase.get(null, theKey, theData, LockMode.DEFAULT)
380                      == OperationStatus.SUCCESS) {
381                  // Recreate the data.
382                  //             DataItem retrievedData = (DataItem) myDataBinding
383                  //                              .entryToObject(theData);
384                  DataItem retrievedData = new DataItem();
385                  retrievedData.setData(new String(theData.getData()));
386                  //add the namespace attribute to DataItem
387                  retrievedData.setNamespace(this.getDatabaseId());
388                  retrievedData.setUserKey(new UserKey(key));
389                  return retrievedData;
390              } else {
391                  return null;
392              }
393          } catch (Exception e) {
394              // Exception handling goes here
395              logger.error("Get: Error hapened: " + e);
396              return null;
397          }
398      }
399
400      public boolean putDataItem(DataItem data) {
401          //extract key from data
402          long key = data.getUserKey().getId();
403          //check, if connection to database is already open
404          if (myDatabase == null) {
405              this.openDatabase();
406          }
407          try {
408              DatabaseEntry theKey = new DatabaseEntry();
409              DatabaseEntry theData = new DatabaseEntry(data.getData().getBytes(
410                      "UTF-8"));
411              //TODO: remove the namespace attribute from data, because it would
412              // need
413              // unneccessary space on disk...
414              //data.removeValue("namespace");
415              //data.removeValue("userkey");
416
417              // TODO: for this test, only store the real data, nothing else into
418              // db! Can change in future
419              DataItem storedDataItem = new DataItem();
420              storedDataItem.setData(data.getData());
421
422              myKeyBinding.objectToEntry(new Long(key), theKey);
423              //             myDataBinding.objectToEntry(storedDataItem, theData);
424              OperationStatus outcome = myDatabase.put(null, theKey, theData);
425              sync();
426              //TODO: revert namespace attribute so we don't remove any data from
427              // original DataItem
428              //data.putValue("namespace", this.getDatabaseId());
429              //data.setUserKey(new UserKey(key));
430              if (outcome == OperationStatus.SUCCESS) {
431                  return true;
432              } else {
433                  logger.error("Put: Error hapened: " + outcome);
434                  return false;
435              }
```

```
436         } catch (Exception e) {
437             // Exception handling goes here
438             logger.error("Put: Error hapened: " + e);
439             return false;
440         }
441     }
442
443     public boolean deleteDataItem(DataItem data) {
444         //extract key from data
445         long key = data.getUserKey().getId();
446         //check, if connection to database is already open
447         if (myDatabase == null) {
448             this.openDatabase();
449         }
450         try {
451             DatabaseEntry theKey = new DatabaseEntry();
452             myKeyBinding.objectToEntry(new Long(key), theKey);
453             // Perform the deletion. All records that use this key are deleted.
454             OperationStatus outcome = myDatabase.delete(null, theKey);
455             sync();
456             if (outcome == OperationStatus.SUCCESS) {
457                 return true;
458             } else {
459                 return false;
460             }
461         } catch (Exception e) {
462             // Exception handling goes here
463             logger.error("Delete: Error hapened: " + e);
464             return false;
465         }
466     }
467
468     public int getDbSize() {
469         // TODO: As of update from je 1.5.1 -> 1.5.3 this does not work
470         // anymore... new ideas needed
471         //         StatsConfig statConfig = new StatsConfig();
472         //         statConfig.setFast(false);
473         //         if (myDatabase == null) {
474         //             this.openDatabase();
475         //         }
476         //         try {
477         //             DatabaseStats dbStat = myDatabase.getStats(statConfig);
478         //             int size = dbStat.;
479         //             logger.info("Current BerkeleyDB size: " + size);
480         //             return size;
481         //         } catch (DatabaseException e) {
482         //             // TODO Auto-generated catch block
483         //             logger.error("get DB Size error: " + e);
484         //         }
485         return -1;
486     }
487
488     private void sync() throws DatabaseException {
489         if (syncCounter > requestsSync) {
490             myDbEnvironment.sync();
491             syncCounter = 0;
492         }
493         syncCounter++;
494     }
495
496     /**
497      * If database is still closed, <code>databaseId</code> is null
498      * <p>
499      * If database is opened and running, <code>databaseId</code> contains an
500      * unique Id (namespace like "ch.nix.accountDB")
501      *
502      * @return Returns the databaseId.
503      */
504     public String getDatabaseId() {
505         return databaseId;
506     }
507
508     public void setDatabaseId(String id) {
509         //TODO check, if id does not contain any reserved id
510         databaseId = id;
511     }
512
513 }
```

---

**File 30:** keypartition.KeyGroup

```
1 /*
2  * Created on Aug 2, 2004 10:46:25 AM
```

```
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 138 $
 8  * Location:    $URL: KeyGroup.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16
17 package ch.nix.ert.keypartition;
18
19 import ch.nix.ert.message.InternalKey;
20 import ch.nix.ert.message.Node;
21
22 /**
23  * @author michi
24  *
25  * TODO Edit this text for a better description
26  */
27 public class KeyGroup implements Comparable {
28
29     private InternalKey keyGroupId = null;
30
31     private Node uplinkNode, downlinkNode;
32
33     /**
34      * @param keyGroupId
35      *              sets the unique ID of this KeyGroup.
36      */
37     public KeyGroup(InternalKey keyGroupId) {
38         this.keyGroupId = keyGroupId;
39     }
40
41     /**
42      * @return Returns the KeyGroup ID.
43      */
44     public InternalKey getKeyGroupId() {
45         return keyGroupId;
46     }
47
48     /**
49      * @return Returns the downlinkNode.
50      */
51     public Node getDownlinkNode() {
52         return downlinkNode;
53     }
54
55     /**
56      * @return Returns the uplinkNode.
57      */
58     public Node getUplinkNode() {
59         return uplinkNode;
60     }
61
62     /* (non-Javadoc)
63      * @see java.lang.Comparable#compareTo(java.lang.Object)
64      */
65     public int compareTo(Object o) {
66         KeyGroup keyGroup = (KeyGroup) o;
67         return getKeyGroupId().compareTo(keyGroup.getKeyGroupId());
68     }
69
70     public String toString() {
71         return getKeyGroupId().toString();
72     }
73 }
```

---

**File 31:** keypartition.KeyPartition

```
1 /*
2  * Created on Aug 2, 2004 12:28:43 PM
3  *
4  * Project:     ExtendendRootTree
5  * Last Change: $Date: 2004-11-12 12:46:25 +0100 (Fri, 12 Nov 2004) $
6  * Changed by:  $Author: michi $
7  * Revision:    $Rev: 173 $
8  * Location:    $URL: KeyPartition.java $
```

```
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16  package ch.nix.ert.keypartition;
17
18  import org.apache.log4j.Logger;
19
20  import ch.nix.ert.base.Dispatcher;
21  import ch.nix.ert.base.ServiceManager;
22  import ch.nix.ert.message.DataItem;
23  import ch.nix.ert.message.InternalKey;
24  import ch.nix.ert.message.Node;
25  import ch.nix.ert.message.UserKey;
26
27  /**
28   * TODO Edit this text for a better description
29   *
30   * @author michi
31   */
32  public class KeyPartition {
33
34      // logger:
35      static Logger logger = Logger.getLogger(KeyPartition.class);
36
37      static private int keyPartitionIdCounter = 0;
38
39      private int keyPartitionId = 0;
40
41      private Node uplinkNode, downlinkNode;
42
43      private KeyGroup firstKeyGroup;
44
45      private KeyGroup lastKeyGroup;
46
47      private Dispatcher dispatcher = ServiceManager.getInstance()
48              .getDispatcher();
49
50      public KeyPartition() {
51          keyPartitionId = keyPartitionIdCounter++;
52      }
53
54      public int getId() {
55          return keyPartitionId;
56      }
57
58      /**
59       * TODO During register process not all KeyGroups are necessary, only the
60       * upper, lower one. We can remove all KeyGroups in between to enhance
61       * performance, so we basically don't need any TreeMap for this class, as we
62       * only need to hold the upper, lower bound? It could be problematic, if we
63       * want to decide, if a special KeyGroup is occupied? (E.g. to replicate it
64       * to another node? If we only now the upper, lower KeyGroups, we can't say
65       * if a dedicated key-group is contained or if it is just empty) <br>
66       * Also if we need to broadcast for all KeyGroups, that we are responsible
67       * for it, could be a bit tricky, if we don't know all key-groups
68       *
69       * @param newKeyGroup
70       */
71      public synchronized void register(KeyGroup newKeyGroup) {
72          boolean hasChanged = add(newKeyGroup);
73          if (hasChanged) {
74              publish();
75          }
76      }
77
78      public synchronized boolean add(KeyGroup newKeyGroup) {
79          //keyGroupList.put(newKeyGroup.getKeyGroupId(), newKeyGroup);
80          boolean hasChanged = false;
81
82          if (firstKeyGroup == null
83                  || newKeyGroup.compareTo(getFirstKeyGroup()) < 0) {
84              firstKeyGroup = newKeyGroup;
85              hasChanged = true;
86          }
87          if (lastKeyGroup == null
88                  || newKeyGroup.compareTo(getLastKeyGroup()) > 0) {
89              lastKeyGroup = newKeyGroup;
90              hasChanged = true;
91          }
92          return hasChanged;
93      }
```

```
 94
 95     public void publish() {
 96         // first publish the information at the keypartitionManager
 97         KeyPartitionManager.getInstance().register(this);
 98
 99         if (getFirstKeyGroup() != null && getLastKeyGroup() != null) {
100             DataItem publish = new DataItem("publish");
101             publish.setKeyPartitionId(keyPartitionId);
102             publish.setFirstInternalKey(getFirstKeyGroup().getKeyGroupId());
103             publish.setLastInternalKey(getLastKeyGroup().getKeyGroupId());
104             dispatcher.addWork(publish);
105         }
106     }
107
108     //    public synchronized void unregister(KeyGroup oldKeyGroup) {
109     //        dispatcher.addWork(new DataItem("unpublish", this));
110     //        keyGroupList.remove(oldKeyGroup.getKeyGroupId());
111     //    }
112
113     public synchronized void unpublish() {
114         DataItem unpublish = new DataItem("unpublish");
115         unpublish.setKeyPartitionId(keyPartitionId);
116         unpublish.setFirstInternalKey(getFirstKeyGroup().getKeyGroupId());
117         unpublish.setLastInternalKey(getLastKeyGroup().getKeyGroupId());
118         dispatcher.addWork(unpublish);
119     }
120
121     public void addKeyPartition(KeyPartition keyPartition) {
122         logger.info("Add: " + keyPartition + "(ID:" + keyPartition.getId()
123                 + ") to me: " + this + "(ID:" + this.getId() + ")");
124         // I only add the first and last KeyGroup
125         boolean hasChanged1 = false;
126         boolean hasChanged2 = false;
127         hasChanged1 = this.add(keyPartition.getFirstKeyGroup());
128         hasChanged2 = this.add(keyPartition.getLastKeyGroup());
129         if (hasChanged1 || hasChanged2) {
130             publish();
131         }
132     }
133
134     /**
135      *
136      *
137      * @param keyPartition
138      * @return the cutoff part of this KeyPartition, null if nothing could be
139      *         cut off
140      */
141     public KeyPartition subtractKeyPartition(KeyPartition keyPartition) {
142         logger.info("Subtract: " + keyPartition + "(ID:" + keyPartition.getId()
143                 + ") from me: " + this + "(ID:" + this.getId() + ")");
144         KeyGroup cutoffFirst, cutoffLast;
145         KeyPartition cutoff = new KeyPartition();
146         KeyGroup firstKeyGroupOfSubstract = keyPartition.getFirstKeyGroup();
147         KeyGroup lastKeyGroupOfSubstract = keyPartition.getLastKeyGroup();
148         // 4 cases
149         if (firstKeyGroupOfSubstract.compareTo(getFirstKeyGroup()) <= 0
150                 && isResponsible(lastKeyGroupOfSubstract)) {
151             cutoffFirst = this.getFirstKeyGroup();
152             cutoffLast = lastKeyGroupOfSubstract;
153
154             cutoff.add(cutoffFirst);
155             cutoff.add(cutoffLast);
156             //register(new KeyGroup(lastKeyGroup.getKeyGroupId()
157             //        .getSubsequentInternalKey()));
158             this.firstKeyGroup = new KeyGroup(lastKeyGroupOfSubstract
159                     .getKeyGroupId().getSubsequentInternalKey());
160             this.publish();
161             return cutoff;
162         }
163
164         if (isResponsible(firstKeyGroupOfSubstract)
165                 && lastKeyGroupOfSubstract.compareTo(getLastKeyGroup()) >= 0) {
166             cutoffFirst = firstKeyGroupOfSubstract;
167             cutoffLast = this.getLastKeyGroup();
168
169             cutoff.add(cutoffFirst);
170             cutoff.add(cutoffLast);
171             //register(new KeyGroup(firstKeyGroup.getKeyGroupId()
172             //        .getPreviousInternalKey()));
173             this.lastKeyGroup = new KeyGroup(firstKeyGroupOfSubstract
174                     .getKeyGroupId().getPreviousInternalKey());
175             this.publish();
176             return cutoff;
177         }
178         if (isResponsible(firstKeyGroupOfSubstract)
```

```
179                   && isResponsible(lastKeyGroupOfSubstract)) {
180              // not supported now... we can't remove from the MIDDLE of the
181              // keygroup as this would imply a split!
182              logger
183                      .warn("Substracting this Keypartition would split it into 2 parts,"
184                          + " this is not supported!");
185              return null;
186          }
187          // this means both borders are outside of this KeyPartition, so we have
188          // to get rid off the whole KeyPartition
189          logger
190                  .info("The whole KeyPartition has been substracted, this is ok, but "
191                      + "could also be a failure?");
192          unpublish();
193          cutoff.addKeyPartition(this);
194          this.firstKeyGroup = null;
195          this.lastKeyGroup = null;
196          return cutoff;
197      }
198
199      //TODO currently returns false, if we could adopt it (adopting is
200      // currently not supported)
201      public boolean isResponsible(DataItem data) {
202          InternalKey internalKey = data.getInternalKey();
203          if (internalKey == null) {
204              UserKey userKey = data.getUserKey();
205              if (userKey == null) {
206                  // we have not found any data to see if this dataItem could fit
207                  // into our Keygroups, so we are not responsible for it
208                  return false;
209              }
210              internalKey = userKey.getInternalKey();
211          }
212          return isResponsible(internalKey);
213      }
214
215      public boolean isResponsible(KeyGroup keyGroup) {
216          return isResponsible(keyGroup.getKeyGroupId());
217      }
218
219      public boolean isResponsible(InternalKey internalKey) {
220          //      check if DataItem falls between my range of keyGroups
221          //check if DataItem is bigger than my lowest value:
222          if (internalKey.compareTo(getFirstKeyGroup().getKeyGroupId()) >= 0) {
223              //check if DataItem is smaller than my highest value:
224              if (internalKey.compareTo(getLastKeyGroup().getKeyGroupId()) <= 0) {
225                  return true;
226              }
227          }
228          // DataItem does not fall into our range, now check if we could adopt
229          // this new DataItem
230          // This does not need to be done (perhaps) as we now ALWAYS extend the
231          // KeyPartition to the highest possible value, so there is never a hole
232          // in the KeyPartition. Example:
233          // If node 1 has keyGroup 1,2,3
234          // If node 2 has keyGroup 100,102,110
235          // Then KeyPartition of node 1 is: [1, 99]
236          // and KeyPartition of node 2 is: [100, 110]
237          // So this "return false" is correct then and no need to further expand
238          // the implementation of this method
239          return false;
240      }
241
242      public KeyGroup getFirstKeyGroup() {
243          return firstKeyGroup;
244      }
245
246      public KeyGroup getLastKeyGroup() {
247          return lastKeyGroup;
248      }
249
250      protected void finalize() throws Throwable {
251          logger.info("Finalized called for KeyPartition ID: " + this.getId());
252          try {
253              this.unpublish(); // close all open published infos on this
254              // object!
255          } finally {
256              super.finalize();
257          }
258      }
259
260      public String toString() {
261          return "ID" + getId() + ":" + getFirstKeyGroup().toString() + "-"
262                  + getLastKeyGroup().toString();
263      }
```

```
264
265 }
```

---

**File 32:** keypartition.KeyPartitionManager

```java
 1 /*
 2  * Created on Aug 2, 2004 10:48:16 AM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-11-12 12:46:25 +0100 (Fri, 12 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 173 $
 8  * Location:    $URL: KeyPartitionManager.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16
17 package ch.nix.ert.keypartition;
18
19 import java.util.Iterator;
20 import java.util.SortedMap;
21 import java.util.TreeMap;
22
23 import org.apache.log4j.Logger;
24
25 import ch.nix.ert.message.DataItem;
26 import ch.nix.ert.message.InternalKey;
27 import ch.nix.ert.message.UserKey;
28
29 /**
30  * TODO Has to check, if this node is responsible for the requested key (by
31  * looking at all his key-partitions) and otherwise has to forward the request
32  * to the corresponding node
33  *
34  * @author michi
35  */
36 public class KeyPartitionManager {
37
38     // logger:
39     static Logger logger = Logger.getLogger(KeyPartitionManager.class);
40
41     private TreeMap keyGroupList = new TreeMap();
42
43     static private KeyPartitionManager instance = null;
44
45     static boolean firstNode = false;
46
47     private int dbSize = 0;
48
49     private int putRequestsPerMinute = 0;
50
51     private int putRequests = 0;
52
53     private int getRequestsPerMinute = 0;
54
55     private int getRequests = 0;
56
57     private int deleteRequestsPerMinute = 0;
58
59     private int deleteRequests = 0;
60
61     protected boolean isJoining = false;
62
63     private KeyPartitionService service;
64
65     private MigrationServer migrationServer;
66
67     private MigrationClient migrationClient;
68
69     private KeyPartitionManager() {
70     }
71
72     static public synchronized KeyPartitionManager getInstance() {
73         if (null == instance) {
74             //System.out.println("Instanciation "+instance+" of
75             // KeyPartitionManager: "+Thread.currentThread().getName());
76             instance = new KeyPartitionManager();
77             instance.service = KeyPartitionService.getInstance();
```

```
78              instance.migrationServer = MigrationServer.getInstance();
79              instance.migrationClient = MigrationClient.getInstance();
80              //System.out.println("Instance set: "+instance+" of
81              // KeyPartitionManager: "+Thread.currentThread().getName());
82              //setting up instance
83              DataItem dbSizeItem = new DataItem();
84              dbSizeItem.setOperation("getsizedb");
85              DataItem callBackItem = new DataItem();
86              callBackItem.setOperation("callcontinous");
87              callBackItem.setDelay(60000);
88              callBackItem.setNested(dbSizeItem);
89
90              DataItem statisticItem = new DataItem();
91              statisticItem.setOperation("updatestatistic");
92              DataItem callBackItem2 = new DataItem();
93              callBackItem2.setOperation("callcontinous");
94              callBackItem2.setDelay(60000);
95              callBackItem2.setNested(statisticItem);
96
97              //registering callback for database size
98              instance.service.transactionBegin(callBackItem);
99              //registering callback for statistic update
100             instance.service.transactionBegin(callBackItem2);
101             //System.out.println("Finished Instanciation "+instance+" of
102             // KeyPartitionManager: "+Thread.currentThread().getName());
103             logger.info("KeyPartitionManager initialized");
104         }
105         return instance;
106     }
107
108     public void register(KeyPartition newKeyPartition) {
109         // first check, if we already have this "new" KeyPartition in our list:
110         for (Iterator iter = keyGroupList.values().iterator(); iter.hasNext();) {
111             KeyPartition keyPartition = (KeyPartition) iter.next();
112             if (keyPartition.getId() == newKeyPartition.getId()) {
113                 iter.remove();
114             }
115         }
116         keyGroupList.put(newKeyPartition.getFirstKeyGroup().getKeyGroupId(),
117                 newKeyPartition);
118     }
119
120     public void unregister(KeyPartition oldKeyPartition) {
121         oldKeyPartition.unpublish();
122         keyGroupList.remove(oldKeyPartition.getFirstKeyGroup().getKeyGroupId());
123     }
124
125     public KeyPartition responsible(DataItem dataItem) {
126         InternalKey internalKey = dataItem.getInternalKey();
127         if (internalKey == null) {
128             UserKey userKey = dataItem.getUserKey();
129             if (userKey == null) {
130                 return null;
131             }
132             return responsible(userKey.getInternalKey());
133         } else {
134             return responsible(internalKey);
135         }
136     }
137
138     public KeyPartition responsible(InternalKey internalKey) {
139         KeyPartition actualKeyPartition = null;
140         SortedMap headMap;
141         if (keyGroupList.isEmpty()) {
142             // this means there is no KeyGroup there, so i return null
143             return null;
144         } else {
145             headMap = keyGroupList.headMap(internalKey
146                     .getSubsequentInternalKey());
147             //if this beast is not empty, get the last key of it:
148             if (!headMap.isEmpty()) {
149                 actualKeyPartition = (KeyPartition) keyGroupList.get(headMap
150                         .lastKey());
151             }
152         }
153         if (actualKeyPartition == null) {
154             return null;
155         }
156         if (actualKeyPartition.isResponsible(internalKey)) {
157             return actualKeyPartition;
158         } else {
159             return null;
160         }
161     }
162
```

```
163      /**
164       *
165       *
166       * @param keyPartition
167       * @return the cutoff part of this KeyPartition, null if nothing could be
168       *         cut off
169       */
170      public KeyPartition substract(KeyPartition keyPartition) {
171          // first find the correct keypartition in my list where I can cutoff
172          // something
173          InternalKey firstKey = keyPartition.getFirstKeyGroup().getKeyGroupId();
174          InternalKey lastKey = keyPartition.getLastKeyGroup().getKeyGroupId();
175          KeyPartition keyPartitionToCut = responsible(firstKey);
176          if (keyPartitionToCut == null) {
177              keyPartitionToCut = responsible(lastKey);
178          }
179          if (keyPartitionToCut == null) {
180              // nothing to do...
181              return null;
182          } else {
183              return keyPartitionToCut.subtractKeyPartition(keyPartition);
184          }
185      }
186
187      public void add(KeyPartition keyPartition) {
188          // I can only add a KeyPartition, if it will extend another one without
189          // any hole in between, as otherwise we would have a splited
190          // KeyPartition.
191
192          if (keyGroupList.isEmpty()) {
193              // this means we have the very first keyPartition here so we should
194              // add it!
195              logger.info("add: First KeyPartition to add: " + keyPartition);
196              register(keyPartition);
197              keyPartition.publish();
198          } else {
199
200              // as adding a new KeyPartition normally does not overlap with the
201              // already existing KeyPartitions we have to add +1 to each edge for
202              // searching the correct keypartition to add this new one
203              InternalKey firstKey = keyPartition.getFirstKeyGroup()
204                      .getKeyGroupId().getPreviousInternalKey();
205              InternalKey lastKey = keyPartition.getLastKeyGroup()
206                      .getKeyGroupId().getSubsequentInternalKey();
207              KeyPartition keyPartitionToAdd = responsible(firstKey);
208              if (keyPartitionToAdd == null) {
209                  keyPartitionToAdd = responsible(lastKey);
210              }
211              if (keyPartitionToAdd == null) {
212                  // nothing to do...
213                  logger
214                          .info("Could not add keyPartition: "
215                                  + keyPartition
216                                  + " to my own: "
217                                  + this
218                                  + ", because they would result in a KeyPartition split...");
219              } else {
220                  keyPartitionToAdd.addKeyPartition(keyPartition);
221              }
222          }
223      }
224
225      public KeyPartition getFirst() {
226          if (keyGroupList.isEmpty()) {
227              return null;
228          }
229          KeyPartition firstKeyPartition = (KeyPartition) keyGroupList
230                  .get(keyGroupList.firstKey());
231          return firstKeyPartition;
232      }
233
234      public KeyPartition getLast() {
235          if (keyGroupList.isEmpty()) {
236              return null;
237          }
238          KeyPartition lastKeyPartition = (KeyPartition) keyGroupList
239                  .get(keyGroupList.lastKey());
240          return lastKeyPartition;
241      }
242
243      public void get(DataItem data) {
244          if (responsible(data) != null) {
245              data.setOperation("getdb");
246              service.transactionContinue(data);
247              getRequests++;
```

```
248            } else {
249                logger
250                        .info("get: I forward this request, as I'm not responsible for: "
251                                + data);
252                data.setRouting("routing");
253                service.transactionContinue(data);
254            }
255        }
256
257    public void put(DataItem data) {
258        if (responsible(data) != null || migrationClient.isResponsible(data)) {
259            // Yes this node is responsible for this data, store it!
260            data.setOperation("putdb");
261            if (migrationServer.isMigrating()
262                    && migrationServer.isResponsible(data)) {
263                // We are the server and are migrating data to somewhere, so I
264                // have to send the put requests also there
265
266                service.transactionSplit(data, 2);
267
268                DataItem data2 = data.getClone();
269                data2.setDirectDestination(migrationServer.getCurrentClient());
270                data2.setRouting("direct");
271                service.transactionSplit(data2, 2);
272
273            } else {
274                service.transactionContinue(data);
275            }
276            putRequests++;
277            // System.out.println("Put Requests received sofar: " +
278            // putRequests);
279        } else {
280            logger
281                    .info("put: I forward this request, as I'm not responsible for: "
282                            + data);
283            data.setRouting("routing");
284            service.transactionContinue(data);
285        }
286    }
287
288    public void delete(DataItem data) {
289        // TODO Perhaps check first if I'm responsible for this DataItem, and
290        // redirect, if necessary?
291        data.setOperation("deletedb");
292        if (migrationServer.isMigrating()) {
293            // We are the server and are migrating data to somewhere, so I
294            // have to send the put requests also there
295
296            service.transactionSplit(data, 2);
297
298            DataItem data2 = data.getClone();
299            data2.setDirectDestination(migrationServer.getCurrentClient());
300            data2.setRouting("direct");
301            service.transactionSplit(data2, 2);
302
303        } else {
304            service.transactionContinue(data);
305        }
306        deleteRequests++;
307    }
308
309    /**
310     * This method should be only called for the first node in the whole
311     * network. All other nodes should use the "join" method to actually join
312     * the ert-network
313     */
314    public void isFirstNodeOnNetwork() {
315        // check that this method is called ONLY once
316        if (firstNode == false) {
317            // I'm the first node, so I'm responsible for the whole
318            logger.info("This node is the first node on the whole network!");
319            KeyGroup firstKeyGroup = new KeyGroup(new InternalKey(
320                    Long.MIN_VALUE));
321            KeyGroup lastKeyGroup = new KeyGroup(
322                    new InternalKey(Long.MAX_VALUE));
323            KeyPartition firstKeyPartition = new KeyPartition();
324            firstKeyPartition.add(firstKeyGroup);
325            firstKeyPartition.add(lastKeyGroup);
326            register(firstKeyPartition);
327            firstKeyPartition.publish();
328            firstNode = true;
329        }
330    }
331
332    public void joinOnNetwork() {
```

```
333          // we have to join on the network, so we have to migrate some data to
334          // us!
335          MigrationClient.getInstance().join();
336      }
337
338      /**
339       * gives back the load of this current node: takes into account the amount
340       * of items stored in db and later also CPU load, throughput of requests,
341       * ...
342       *
343       * @return value between 0-1 (0 = idle, 1 = on system limit, >1 = heavily
344       *         overloaded)
345       */
346      public double nodeLoad() {
347          // this is hardcoded, so an amount of 100000 stored attributes are
348          // looked as too much
349          //        int maxStoredValues = 100000;
350          //        float store = (float) dbSize / (float) maxStoredValues;
351          //
352          //        //hardcoded put value 500 requests per minute considered as very
353          // high:
354          //        int maxPutPerMinute = 500;
355          //        float putdb = (float) putRequestsPerMinute / (float) maxPutPerMinute;
356          //
357          //        //hardcoded put value 500 requests per minute considered as very
358          // high:
359          //        int maxGetPerMinute = 500;
360          //        float getdb = (float) getRequestsPerMinute / (float) maxGetPerMinute;
361          //
362          //        //hardcoded put value 500 requests per minute considered as very
363          // high:
364          //        int maxDeletePerMinute = 500;
365          //        float deletedb = (float) deleteRequestsPerMinute
366          //                / (float) maxDeletePerMinute;
367
368          // now do a weighting of the values:
369          // store: 0.4
370          // put : 0.2
371          // get : 0.3
372          // delete: 0.1
373          //        System.out.println("DB-Size:" + dbSize + " Puts/min:"
374          //                + putRequestsPerMinute + " Gets/min:" + getRequestsPerMinute
375          //                + " Dels/min:" + deleteRequestsPerMinute);
376          //float load = (float) (store * 0.4 + putdb * 0.2 + getdb * 0.3 +
377          // deletedb * 0.1);
378          //TODO revert to line above
379          //currently for simplicity only take into account the store:
380          //float load = (float) store;
381
382          // For simplicity I change nodeload to be dependant of the number of
383          // KeyGroups: so the load is just the % of the whole key-range:
384          long numberOfKeyGroups = 0;
385          for (Iterator iter = keyGroupList.values().iterator(); iter.hasNext();) {
386              KeyPartition keyPartition = (KeyPartition) iter.next();
387              numberOfKeyGroups += keyPartition.getFirstKeyGroup()
388                      .getKeyGroupId().amountBetween(
389                              keyPartition.getLastKeyGroup().getKeyGroupId());
390          }
391          long maxPossibleKeyRange = InternalKey.getHigherLimit()
392                  - InternalKey.getLowerLimit();
393
394          double load = (double) numberOfKeyGroups / (double) maxPossibleKeyRange;
395          //load = (float) 0.4;
396          return load;
397      }
398
399      /**
400       * Calculates the KeyRange to migrate
401       *
402       * @param remoteLoad
403       *            the load of the other node, which wants to have some data from
404       *            this node
405       * @param uplinkNode
406       *            if true, the node who requests the migration is the
407       *            uplinkNode, if false it is the downlonkNode
408       * @return The DataItem with parameter "migrationCutPointLow" and
409       *         "migrationCutPointHigh" of the keygroups to migrate
410       */
411      public DataItem calculateMigrationData(double remoteLoad, boolean uplinkNode) {
412          double localLoad = nodeLoad();
413          long storedDataItems = (long) (localLoad * 100000);
414
415          // first check how many dataItems we potentionally could migrate:
416          double difference = localLoad - remoteLoad;
417
```

```
418          if (storedDataItems == 0) {
419              //means we have empty node, so we can easily share half of our keys
420              // to the other node
421              return null;
422          }
423
424          // 40% of the load is related to number of dataitems (see nodeload
425          // method)
426          long numberOfDataItems = (long) (difference * 100000 / 2);
427
428          // now calculate the keygroups we could migrate
429          // First calculate the amount of keygroups we currently hold
430          long numberOfKeyGroups = 0;
431          for (Iterator iter = keyGroupList.values().iterator(); iter.hasNext();) {
432              KeyPartition keyPartition = (KeyPartition) iter.next();
433              numberOfKeyGroups += keyPartition.getFirstKeyGroup()
434                      .getKeyGroupId().amountBetween(
435                              keyPartition.getLastKeyGroup().getKeyGroupId());
436          }
437
438          // now calculate the amount of keygroups to migrate:
439          long keygroupsToMigrate = (long) (((double) numberOfKeyGroups /
440                  (double) storedDataItems) *
441                  (double) numberOfDataItems);
442
443          // now calculate the correct keygroup where we should do the cut - we
444          // migrate max one full keypartition, even if keygroupsToMigrate is
445          // larger than the last keypartition
446          DataItem response = new DataItem();
447          if (uplinkNode) {
448              // if the other node is the uplink node, we have to migrate data
449              // from the higher end of the last keyPartition
450              KeyPartition lastKeyPartition = (KeyPartition) keyGroupList
451                      .get(keyGroupList.lastKey());
452
453              InternalKey lastInternalKey = lastKeyPartition.getLastKeyGroup()
454                      .getKeyGroupId();
455
456              InternalKey migrationCutPoint = lastInternalKey
457                      .getNthPreviousInternalKey(keygroupsToMigrate);
458              // now check if this cutPoint is within the lastKeyPartition
459              if (!lastKeyPartition.isResponsible(migrationCutPoint)) {
460                  //ops lastKeyPartition is not responsible for this key, so I
461                  // set
462                  // the migrationCutPoint to the first key!
463                  migrationCutPoint = lastKeyPartition.getFirstKeyGroup()
464                          .getKeyGroupId();
465              }
466              response.setLastInternalKey(lastInternalKey);
467              response.setFirstInternalKey(migrationCutPoint);
468          } else {
469              // if the other node is the downlink node, we have to migrate data
470              // from the lower end of the first keyPartition
471              KeyPartition firstKeyPartition = (KeyPartition) keyGroupList
472                      .get(keyGroupList.firstKey());
473
474              InternalKey firstInternalKey = firstKeyPartition.getFirstKeyGroup()
475                      .getKeyGroupId();
476
477              InternalKey migrationCutPoint = firstInternalKey
478                      .getNthSubsequentInternalKey(keygroupsToMigrate - 1);
479              // now check if this cutPoint is within the lastKeyPartition
480              if (!firstKeyPartition.isResponsible(migrationCutPoint)) {
481                  //ops lastKeyPartition is not responsible for this key, so I
482                  // set
483                  // the migrationCutPoint to the first key!
484                  migrationCutPoint = firstKeyPartition.getLastKeyGroup()
485                          .getKeyGroupId();
486              }
487              response.setFirstInternalKey(firstInternalKey);
488              response.setLastInternalKey(migrationCutPoint);
489          }
490          return response;
491      }
492
493      protected void updatedbSize(int size) {
494          //        System.out.println("Update DB size called! having done: " + size
495          //                + " entries");
496          dbSize = size;
497      }
498
499      protected void updateStatistic() {
500          //        System.out.println("Update Statistic called! having done: "
501          //                + putRequests + " putRequests");
502          putRequestsPerMinute = putRequests;
```

```
503          getRequestsPerMinute = getRequests;
504          deleteRequestsPerMinute = deleteRequests;
505          putRequests = 0;
506          getRequests = 0;
507          deleteRequests = 0;
508      }
509
510      public String toString() {
511          String representation = "( ";
512          for (Iterator iter = keyGroupList.values().iterator(); iter.hasNext();) {
513              KeyPartition keyPartition = (KeyPartition) iter.next();
514              representation = representation + keyPartition + " ";
515          }
516          return representation + ")";
517      }
518
519      //    public boolean isResponsible(DatItem data) {
520      //
521      //    }
522
523 }
```

---

**File 33:** keypartition.KeyPartitionService

```
 1 /*
 2  * Created on Aug 27, 2004 12:29:27 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-11-09 01:55:31 +0100 (Tue, 09 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 165 $
 8  * Location:    $URL: KeyPartitionService.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.keypartition;
17
18 import org.apache.log4j.Logger;
19
20 import ch.nix.ert.base.ServiceBase;
21 import ch.nix.ert.message.DataItem;
22
23 /**
24  * TODO Edit this text for a better description
25  *
26  * @author michi
27  */
28 public class KeyPartitionService extends ServiceBase {
29
30      // logger:
31      static Logger logger = Logger.getLogger(KeyPartitionService.class);
32
33      static private KeyPartitionService instance = null;
34
35      private KeyPartitionManager keyPartitionManager;
36
37      private MigrationClient migrationClient;
38
39      private MigrationServer migrationServer;
40
41      // private TreeMap timeSortedList = new TreeMap();
42
43      private KeyPartitionService() {
44      }
45
46      static public synchronized KeyPartitionService getInstance() {
47          if (null == instance) {
48              instance = new KeyPartitionService();
49              instance.keyPartitionManager = KeyPartitionManager.getInstance();
50              instance.migrationClient = MigrationClient.getInstance();
51              instance.migrationServer = MigrationServer.getInstance();
52              logger.info("KeyPartitionService started");
53          }
54          return instance;
55      }
56
57      /*
58       * (non-Javadoc)
```

```
59        *
60        * @see ch.nix.ert.ServiceInterface#registerAtDispatcher()
61        */
62       public void registerAtDispatcher() {
63           // TODO Auto-generated method stub
64           registerService("get");
65           registerService("put");
66           registerService("delete");
67
68           registerService("updatedbsize");
69           registerService("updatestatistic");
70           registerService("nodeload");
71
72           registerService("firstnode");
73           registerService("join");
74
75           registerService("migrationprobe");
76           //registerService("migrationservercancel");
77           //registerService("migrationclientcancel");
78           registerService("migrationstart");
79           registerService("migrationend");
80           registerService("getmigration");
81           registerService("deletemigration");
82       }
83
84       /*
85        * (non-Javadoc)
86        *
87        * @see java.lang.Runnable#run()
88        */
89       public void run() {
90           //System.out.println("DataStoreService started");
91           Thread.currentThread().setName("ERT.KeyPartitionService");
92           while (true) {
93               DataItem dataItem = null;
94               dataItem = getWork();
95               // logger.debug("Got Work: "+dataItem);
96               String operation = dataItem.getOperationType();
97               //System.out.println("Operation: "+operation+" equals? "+
98               // (operation.equals("nodeload")));
99               if (operation.equals("get")) {
100                  // logger.debug("get: " + dataItem);
101                  keyPartitionManager.get(dataItem);
102              } else if (operation.equals("put")) {
103                  // logger.debug("put: " + dataItem);
104                  //System.out.println("put in keyservice...");
105                  keyPartitionManager.put(dataItem);
106              } else if (operation.equals("delete")) {
107                  // logger.debug("delete: " + dataItem);
108                  keyPartitionManager.delete(dataItem);
109              } else if (operation.equals("getsizedb-end")) {
110                  // logger.debug("getSizeDBResponse: " + dataItem);
111                  keyPartitionManager.updatedbSize(dataItem.getDbSize());
112              } else if (operation.equals("updatestatistic")) {
113                  // logger.debug("udpadateStatistic: " + dataItem);
114                  //  System.out
115                  //   .println("KeyPartitionService:UPDATESTATISTIC called");
116                  keyPartitionManager.updateStatistic();
117              } else if (operation.equals("firstnode")) {
118                  keyPartitionManager.isFirstNodeOnNetwork();
119              } else if (operation.equals("join")) {
120                  keyPartitionManager.joinOnNetwork();
121              } else if (operation.equals("nodeload")) {
122                  // logger.debug("nodeLoad: " + dataItem);
123                  //Node originator = dataItem.getDirectOriginator();
124
125                  double size = keyPartitionManager.nodeLoad();
126                  DataItem response = new DataItem();
127                  response.setNodeLoad(size);
128                  //response.setDirectDestination(originator);
129                  //response.setOperation("nodeloadresponse");
130                  //System.out.println("KeyPartitionService response: " +
131                  // response);
132                  transactionEnd(response);
133                  // logger.debug("nodeLoadResponse: " + response);
134              } else if (operation.equals("nodeload-end")) {
135                  // logger.debug("nodeLoad-end: " + dataItem);
136                  if (keyPartitionManager.isJoining) {
137                      migrationClient.joinOnNetwork(dataItem);
138                  } else {
139                      migrationClient.nodeLoadResponse(dataItem);
140                  }
141              } else if (operation.equals("migrationprobe")) {
142                  // logger.debug("migrationProbe: " + dataItem);
143                  migrationClient.probeServer();
```

```
144 //          } else if (operation.equals("migrationservercancel")) {
145 //              logger.debug("migrationServerCancel: " + dataItem);
146 //              migrationServer.cancelMigration();
147 //          } else if (operation.equals("migrationclientcancel")) {
148 //              logger.debug("migrationClientCancel: " + dataItem);
149 //              migrationClient.cancelMigration();
150             } else if (operation.equals("migrationstart")) {
151                 // logger.debug("migrationStart: " + dataItem);
152                 migrationServer.migrationStart(dataItem);
153             } else if (operation.equals("migrationstart-end")) {
154                 // logger.debug("migrationStartResponse: " + dataItem);
155                 migrationClient.migrationStartResponse(dataItem);
156             } else if (operation.equals("migrationend")) {
157                 // logger.debug("migrationEnd: " + dataItem);
158                 migrationServer.migrationEnd(dataItem);
159             } else if (operation.equals("migrationend-end")) {
160                 // logger.debug("migrationEndResponse: " + dataItem);
161                 migrationClient.migrationEndResponse(dataItem);
162             } else if (operation.equals("getmigration")) {
163                 // logger.debug("getMigration: " + dataItem);
164                 migrationServer.sendMigratedData(dataItem);
165             } else if (operation.equals("getmigration-end")) {
166                 // logger.debug("getMigrationResponse: " + dataItem);
167                 migrationClient.receiveMigratedData(dataItem);
168             } else {
169
170                 logger
171                     .warn("DataItem dropped, because no processing was found for: "
172                         + dataItem);
173             }
174         }
175
176     }
177 }
```

---

### File 34: keypartition.MigrationClient

```
1 /*
2  * Created on Sep 6, 2004 2:59:21 PM
3  *
4  * Project:     ExtendendRootTree
5  * Last Change: $Date: 2004-11-12 12:46:25 +0100 (Fri, 12 Nov 2004) $
6  * Changed by:  $Author: michi $
7  * Revision:    $Rev: 173 $
8  * Location:    $URL: MigrationClient.java $
9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.keypartition;
17
18 import java.util.Calendar;
19
20 import org.apache.log4j.Logger;
21
22 import ch.nix.ert.message.DataItem;
23 import ch.nix.ert.message.InternalKey;
24 import ch.nix.ert.message.Node;
25 import ch.nix.ert.message.UserKey;
26
27 /**
28  * TODO Edit this text for a better description
29  *
30  * @author michi
31  */
32 public class MigrationClient {
33
34     // logger:
35     static Logger logger = Logger.getLogger(MigrationClient.class);
36
37     static private MigrationClient instance = null;
38
39     private boolean joinExecuted = false;
40
41     private MigrationServer migrationServer = null;
42
43     private KeyPartitionService service;
44
45     private boolean isMigratingAsClient = false;
```

```
46
47      private Node migrationNode;
48
49      private KeyPartition migrationKeyPartition;
50
51      private KeyPartitionManager keyPartitionManager;
52
53      private int joinOnNetworkCount = 10;
54
55      private Node highestLoadNode = new Node();
56
57      private double highestLoad = 0;
58
59      private long migrationTime = 0;
60
61      // hardcoded set: timeout migration after 60 sec of inactivity
62      private long migrationTimeout = 60000;
63
64      private MigrationClient() {
65      }
66
67      static public synchronized MigrationClient getInstance() {
68          if (null == instance) {
69              instance = new MigrationClient();
70              instance.service = KeyPartitionService.getInstance();
71              instance.keyPartitionManager = KeyPartitionManager.getInstance();
72              instance.migrationServer = MigrationServer.getInstance();
73              // calling probeServer to get the migration started!
74              //instance.probeServer(new Node("10.0.1.250", 1099));
75              // now start the automatic probe: every 60 sec
76
77              logger.info("MigrationClient initialised");
78          }
79          return instance;
80      }
81
82      public boolean isMigrating() {
83          migrationTimeout();
84          return isMigratingAsClient;
85      }
86
87      public void join() {
88          if (!joinExecuted) {
89              DataItem migrationProbe = new DataItem();
90              migrationProbe.setOperation("migrationprobe");
91              DataItem callBackItem = new DataItem();
92              callBackItem.setOperation("callcontinous");
93              callBackItem.setDelay(20000);
94              callBackItem.setNested(migrationProbe);
95              service.transactionBegin(callBackItem);
96              joinExecuted = true;
97          }
98          joinOnNetwork(null);
99      }
100
101     public void joinOnNetwork(DataItem dataItem) {
102         logger.info("Client: Probe a random Server for Joining Response: "
103                 + dataItem);
104         if (joinOnNetworkCount > 0) {
105             keyPartitionManager.isJoining = true;
106         } else {
107             keyPartitionManager.isJoining = false;
108             initiateMigration(highestLoadNode);
109         }
110         if (dataItem != null) {
111             if (dataItem.getNodeLoad() > highestLoad) {
112                 highestLoadNode = dataItem.getDirectOriginator();
113                 highestLoad = dataItem.getNodeLoad();
114             }
115         }
116         if (joinOnNetworkCount > 0) {
117             // now send to another random node:
118             DataItem sendData = new DataItem("nodeload");
119             sendData.setRouting("routing");
120             sendData.setInternalKey(InternalKey.getRandom());
121             service.transactionBegin(sendData);
122             logger.info("Probe a random Server for Joining: " + sendData);
123             joinOnNetworkCount-;
124         }
125     }
126
127     // TODO must be: probeServer(Node node)
128     public void probeServer() {
129
130         logger.info("Client: Send Probe to Servers. My current KeyPartitions: "
```

```
131                   + keyPartitionManager);
132           // first check for node-load of my Server
133           KeyPartition firstKeyPartition = keyPartitionManager.getFirst();
134           KeyPartition lastKeyPartition = keyPartitionManager.getLast();
135
136           if (firstKeyPartition != null) {
137               InternalKey downlink = firstKeyPartition.getFirstKeyGroup()
138                       .getKeyGroupId().getPreviousInternalKey();
139               if (keyPartitionManager.responsible(downlink) == null) {
140                   // avoid sending messages to myself, if I'm on the
141                   // lowest/highest edge of internal-Keys
142                   DataItem sendDataDownlink = new DataItem("nodeload");
143                   sendDataDownlink.setRouting("routing");
144                   sendDataDownlink.setInternalKey(downlink);
145                   logger.info("Client: Probe Server: " + sendDataDownlink);
146                   service.transactionBegin(sendDataDownlink);
147               }
148           }
149           if (lastKeyPartition != null) {
150               InternalKey uplink = lastKeyPartition.getLastKeyGroup()
151                       .getKeyGroupId().getSubsequentInternalKey();
152               if (keyPartitionManager.responsible(uplink) == null) {
153                   // avoid sending messages to myself, if I'm on the
154                   // lowest/highest edge of internal-Keys
155                   DataItem sendDataUpLink = new DataItem("nodeload");
156                   sendDataUpLink.setRouting("routing");
157                   sendDataUpLink.setInternalKey(uplink);
158                   logger.info("Client: Probe Server: " + sendDataUpLink);
159                   service.transactionBegin(sendDataUpLink);
160               }
161           }
162           if (firstKeyPartition == null && lastKeyPartition == null) {
163               // means we have not yet found a node to join, so we have to repeat
164               // the node-join algorithm:
165               //
166               // first reset the join-related variables:
167               joinOnNetworkCount = 10;
168               highestLoad = 0;
169               join();
170           }
171           //System.out.println("Data sent");
172       }
173
174       public void nodeLoadResponse(DataItem dataItem) {
175           logger.info("Client: Probe Server Response: " + dataItem);
176           if (!isMigrating()) {
177               // we got an nodeLoadResponse, now check, if we can migrate...
178               double localLoad = keyPartitionManager.nodeLoad();
179               if (localLoad + 0.02 < dataItem.getNodeLoad() || (localLoad == 0)) {
180                   // the Server has at least a 5% higher load than we have, or the
181                   // other node is almost idle, so we can easily so we should
182                   // migrate
183                   initiateMigration(dataItem.getDirectOriginator());
184                   //instance.probeServer(new Node("192.168.0.111", 1099));
185               }
186           }
187       }
188
189       public void initiateMigration(Node node) {
190           if (!migrationServer.isMigrating()) {
191               migrationNode = node;
192               DataItem sendData = new DataItem("migrationstart");
193               sendData.setDirectDestination(migrationNode);
194               KeyPartition firstKeyPartition = keyPartitionManager.getFirst();
195               if (firstKeyPartition != null) {
196                   // this has to be done, so the Server Node finds out if this
197                   // node is
198                   // the uplink, or downlink node
199                   sendData.setInternalKey(firstKeyPartition.getFirstKeyGroup()
200                           .getKeyGroupId());
201               }
202               sendData.setRouting("direct");
203               sendData.setNodeLoad(keyPartitionManager.nodeLoad());
204               logger.info("Client: Migration-Request Initiated: " + sendData);
205               service.transactionBegin(sendData);
206           }
207       }
208
209       public void migrationStartResponse(DataItem dataItem) {
210           if (isMigratingAsClient == false && !migrationTimeout()) {
211               logger.info("Client: Migration Start Response: " + dataItem);
212               isMigratingAsClient = true;
213               // first check, if we want to migrate this amount of data?
214               // TODO: check - currently we always accept.
215
```

```
216                 // we first upload the the migration-data, and if successfully
217                 // finished will publish the new information!
218                 //
219                 // first we generate a temporary KeyPartition (which is not
220                 // published yet) and use it to fill up the data:
221                 migrationKeyPartition = new KeyPartition();
222                 // add the proposed data-Range of the data to migrate:
223                 migrationKeyPartition.add(new KeyGroup(dataItem
224                         .getFirstInternalKey()));
225                 migrationKeyPartition.add(new KeyGroup(dataItem
226                         .getLastInternalKey()));
227
228                 // set the DataItem to the highest value:
229                 DataItem getData = new DataItem("getmigration");
230                 getData.setUserKey(dataItem.getLastInternalKey().getLastUserKey()
231                         .getSubsequentUserKey());
232                 getData.setRouting("direct");
233                 getData.setDirectDestination(migrationNode);
234                 service.transactionBegin(getData);
235             }
236         }
237
238     private boolean migrationTimeout() {
239         long currentTime = Calendar.getInstance().getTimeInMillis();
240         if (currentTime > migrationTime + migrationTimeout
241                 && migrationTime != 0) {
242             logger
243                     .warn("Client: Migration has been canceled, because of an time-out "
244                         + "during transfering of DataItems");
245             isMigratingAsClient = false;
246             migrationKeyPartition = null;
247             migrationTime = 0;
248             return true;
249         }
250         return false;
251     }
252
253     public void receiveMigratedData(DataItem dataItem) {
254         if (isMigratingAsClient) {
255             migrationTime = Calendar.getInstance().getTimeInMillis();
256             // first check if migrated dataItem is within our
257             // MigrationKeyPartition
258             if (migrationKeyPartition.isResponsible(dataItem)) {
259                 // this dataItem is part of this Migration and we store it into
260                 // db
261                 UserKey currentUserKey = dataItem.getUserKey();
262                 dataItem.setOperation("putdbnoack");
263                 service.transactionBegin(dataItem);
264
265                 // now request the next dataItem
266                 DataItem getData = new DataItem("getmigration");
267                 getData.setUserKey(currentUserKey);
268                 getData.setDirectDestination(migrationNode);
269                 getData.setRouting("direct");
270                 service.transactionBegin(getData);
271             } else {
272                 // means we are out of range -> means we have finished the
273                 // migration and could commit it
274                 logger.info("Client: Migration finished");
275                 DataItem sendData = new DataItem("migrationend");
276                 sendData.setDirectDestination(migrationNode);
277                 sendData.setRouting("direct");
278                 service.transactionBegin(sendData);
279             }
280         }
281     }
282
283     public boolean isResponsible(DataItem dataItem) {
284         if (migrationKeyPartition == null) {
285             return false;
286         }
287         return migrationKeyPartition.isResponsible(dataItem);
288     }
289
290     public void migrationEndResponse(DataItem dataItem) {
291         if (isMigratingAsClient) {
292             logger.info("Client: Migration End Response: " + dataItem);
293             keyPartitionManager.add(migrationKeyPartition);
294             migrationKeyPartition = null;
295             migrationNode = null;
296             isMigratingAsClient = false;
297             migrationTime = 0;
298         }
299     }
300
```

```
301 }
```

---

**File 35:** keypartition.MigrationServer

```
 1 /*
 2  * Created on Sep 6, 2004 2:59:39 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-11-12 12:46:25 +0100 (Fri, 12 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 173 $
 8  * Location:    $URL: MigrationServer.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.keypartition;
17
18 import java.util.Calendar;
19
20 import org.apache.log4j.Logger;
21
22 import ch.nix.ert.message.DataItem;
23 import ch.nix.ert.message.InternalKey;
24 import ch.nix.ert.message.Node;
25
26 /**
27  * TODO Edit this text for a better description
28  *
29  * @author michi
30  */
31 public class MigrationServer {
32
33     // logger:
34     static Logger logger = Logger.getLogger(MigrationServer.class);
35
36     static private MigrationServer instance = null;
37
38     private MigrationClient migrationClient = null;
39
40     private Node currentClient = null;
41
42     private KeyPartitionService service;
43
44     private boolean requestMigration = false;
45
46     private boolean isMigratingAsServer = false;
47
48     private KeyPartitionManager keyPartitionManager;
49
50     private KeyPartition migrationKeyPartition;
51
52     private long requestMigrationTime = 0;
53
54     private long migrationTime = 0;
55
56     // hardcoded set: timeout migration after 60 sec of inactivity
57     private long requestMigrationTimeout = 60000;
58
59     private long migrationTimeout = 60000;
60
61     private MigrationServer() {
62     }
63
64     static public synchronized MigrationServer getInstance() {
65         if (null == instance) {
66             instance = new MigrationServer();
67             instance.service = KeyPartitionService.getInstance();
68             instance.keyPartitionManager = KeyPartitionManager.getInstance();
69             instance.migrationClient = MigrationClient.getInstance();
70             logger.info("MigrationServer initialised");
71         }
72         return instance;
73     }
74
75     public boolean isMigrating() {
76         requestMigrationTimeout();
77         migrationTimeout();
78         return requestMigration;
```

```
79        }
80
81        public Node getCurrentClient() {
82            return currentClient;
83        }
84
85        private boolean requestMigrationTimeout() {
86            long currentTime = Calendar.getInstance().getTimeInMillis();
87            if (currentTime > requestMigrationTime + requestMigrationTimeout
88                    && requestMigrationTime != 0) {
89                if (requestMigration && !isMigratingAsServer) {
90                    logger
91                            .info("Server: Migration has been canceled, because of an time-out "
92                                    + "while waiting for the first DataItem to be transfered");
93                    requestMigration = false;
94                    isMigratingAsServer = false;
95                    migrationKeyPartition = null;
96                    currentClient = null;
97                    requestMigrationTime = 0;
98                    migrationTime = 0;
99                    return true;
100               }
101           }
102           return false;
103       }
104
105       private boolean migrationTimeout() {
106           long currentTime = Calendar.getInstance().getTimeInMillis();
107           if (currentTime > migrationTime + migrationTimeout
108                   && migrationTime != 0) {
109               logger
110                       .info("Server: Migration has been canceled, because of an time-out "
111                               + "during transfering of DataItems");
112               requestMigration = false;
113               isMigratingAsServer = false;
114               migrationKeyPartition = null;
115               currentClient = null;
116               requestMigrationTime = 0;
117               migrationTime = 0;
118               return true;
119           }
120           return false;
121       }
122
123       public void migrationStart(DataItem dataItem) {
124           if (requestMigration == false && !migrationClient.isMigrating()
125                   && !requestMigrationTimeout() && !migrationTimeout()) {
126               requestMigrationTime = Calendar.getInstance().getTimeInMillis();
127               logger.info("Server: Migration Start: " + dataItem);
128               // now check which data we can can migrate
129               //
130               // first check if the node who requests the data is uplink or
131               // downlink node
132               InternalKey remoteInternalKey = dataItem.getInternalKey();
133               boolean uplink = true;
134               if (remoteInternalKey == null) {
135                   // no InternalKey set, so remote does not yet know so we have to
136                   // decide :-)
137                   uplink = true;
138               } else if (dataItem.getInternalKey().compareTo(
139                       keyPartitionManager.getFirst().getFirstKeyGroup()
140                           .getKeyGroupId()) < 0) {
141                   uplink = false;
142               }
143
144               DataItem remoteData = keyPartitionManager.calculateMigrationData(
145                       dataItem.getNodeLoad(), uplink);
146               dataItem.setLastInternalKey(remoteData.getLastInternalKey());
147               dataItem.setFirstInternalKey(remoteData.getFirstInternalKey());
148
149               migrationKeyPartition = new KeyPartition();
150               // add the proposed data-Range of the data to migrate:
151               migrationKeyPartition.add(new KeyGroup(dataItem
152                       .getFirstInternalKey()));
153               migrationKeyPartition.add(new KeyGroup(dataItem
154                       .getLastInternalKey()));
155               dataItem.removeNodeLoad();
156               service.transactionEnd(dataItem);
157               currentClient = dataItem.getTransactionNodeId();
158               requestMigration = true;
159           }
160       }
161
162   //    public void cancelMigration() {
163   //        if (requestMigration && !isMigratingAsServer) {
```

```
164     //              logger
165     //                      .info("Server: Migration has been canceled, because of an time-out");
166     //              requestMigration = false;
167     //              isMigratingAsServer = false;
168     //              migrationKeyPartition = null;
169     //          }
170     //      }
171
172     public void sendMigratedData(DataItem dataItem) {
173         if (requestMigration) {
174             migrationTime = Calendar.getInstance().getTimeInMillis();
175             isMigratingAsServer = true;
176             //logger.info("sendMigratedData: " + dataItem);
177             //TODO: as getprev will ALWAYS omit the first key, we have to
178             // check, that for first execution we actually do a normal
179             // get-Operation...
180             // Or because we currently let the whole migration finish untill we
181             // change the Responsibilities of a KeyRange, we can always fetch
182             // backwards...
183
184             // forward this request to DataStore layer, so he actually returns
185             // the data to the Migration-Client
186             dataItem.setOperation("getprevdb");
187             service.transactionContinue(dataItem);
188         }
189     }
190
191     public boolean isResponsible(DataItem dataItem) {
192         if (migrationKeyPartition == null) {
193             return false;
194         }
195         return migrationKeyPartition.isResponsible(dataItem);
196     }
197
198     public void migrationEnd(DataItem dataItem) {
199         if (requestMigration) {
200             logger.info("Server: Migration End: " + dataItem);
201             // now change my KeyPartition
202             keyPartitionManager.substract(migrationKeyPartition);
203
204             // now delete all my obsolete data
205             // TODO... currently no data is deleted, perhaps we can do a general
206             // "cleanup" Service which will delete unused data from db during
207             // low CPU time...
208
209             service.transactionEnd(dataItem);
210             migrationKeyPartition = null;
211             requestMigration = false;
212             isMigratingAsServer = false;
213             requestMigrationTime = 0;
214             migrationTime = 0;
215         }
216     }
217
218 }
```

### File 36: main.MainService

```
 1 /*
 2  * Created on Oct 20, 2004 4:06:01 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-11-12 12:46:25 +0100 (Fri, 12 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 173 $
 8  * Location:    $URL: MainService.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.main;
17
18 import java.awt.Frame;
19 import java.util.ArrayList;
20 import java.util.Arrays;
21
22 import org.apache.log4j.Logger;
23
24 import ch.nix.ert.base.RootLogger;
```

```java
25 import ch.nix.ert.base.ServiceBase;
26 import ch.nix.ert.base.ServiceManager;
27 import ch.nix.ert.message.DataItem;
28
29 /**
30  * MainService class which is executed at startup. <br>
31  * It will initialize the whole system and starts it up
32  *
33  * @author michi
34  */
35 public class MainService extends ServiceBase {
36
37     // logger:
38     static Logger logger = Logger.getLogger(MainService.class);
39
40     static private MainService instance = null;
41
42     static ServiceManager serviceManager;
43
44     static private RootLogger ertLogger;
45
46     static public synchronized MainService getInstance() {
47         if (null == instance) {
48             instance = new MainService();
49             logger.info("MainService started");
50         }
51         return instance;
52     }
53
54     public static void main(String[] args) {
55         // Disable IPv6, as it is not used at the moment
56         System.setProperty("java.net.preferIPv6Addresses", "false");
57         System.setProperty("java.net.preferIPv4Stack", "true");
58         // A possible workaround for an apple bug in kqueue, so I now use the
59         // older select
60         System.setProperty("java.nio.preferSelect", "true");
61
62         ertLogger = new RootLogger();
63         serviceManager = ServiceManager.getInstance();
64
65         MainService mainService = MainService.getInstance();
66         ArrayList parameters = new ArrayList(Arrays.asList(args));
67         boolean success = false;
68
69         if (!parameters.isEmpty() && parameters.contains("-v")) {
70             String version = "$AntReplaceVersion$";
71             String date = "$AntReplaceDate$";
72             System.out.println();
73             System.out.println("ert version:");
74             System.out.println("   Revision: " + version);
75             System.out.println("       Date: " + date);
76             System.out.println();
77             System.exit(0);
78         }
79
80         if (!parameters.isEmpty() && parameters.contains("-c")) {
81             // this means a command has been provided, now get the command
82             if (parameters.size() > 1) {
83                 String command = (String) parameters.get(parameters
84                         .indexOf("-c") + 1);
85                 command = command.toLowerCase();
86                 // only commands allowed are: "join" "first"
87                 if (command.equals("join")) {
88                     logger.info("'-c join' command supplied at startup");
89                     mainService.addWork(new DataItem("join"));
90                     success = true;
91                 }
92                 if (command.equals("first")) {
93                     logger.info("'-c first' command supplied at startup");
94                     mainService.addWork(new DataItem("firstnode"));
95                     success = true;
96                 }
97                 if (command.equals("writebench")) {
98                     logger.info("'-c writebench' command supplied at startup");
99                     String sendStatistics = (String) parameters.get(parameters
100                            .indexOf("-s") + 1);
101                     String dataSize = (String) parameters.get(parameters
102                            .indexOf("-d") + 1);
103                     if (sendStatistics != null) {
104                         DataItem data = new DataItem("writebench");
105                         data.setData(sendStatistics);
106                         data.setDbSize(Integer.parseInt(dataSize));
107                         mainService.addWork(data);
108                         success = true;
109                     }
```

```
110                      }
111                      if (command.equals("readbench")) {
112                          logger.info("'-c readbench' command supplied at startup");
113                          String sendStatistics = (String) parameters.get(parameters
114                                  .indexOf("-s") + 1);
115                          if (sendStatistics != null) {
116                              DataItem data = new DataItem("readbench");
117                              data.setData(sendStatistics);
118                              mainService.addWork(data);
119                              success = true;
120                          }
121                      }
122                      if (command.equals("deletebench")) {
123                          logger.info("'-c deletebench' command supplied at startup");
124                          String sendStatistics = (String) parameters.get(parameters
125                                  .indexOf("-s") + 1);
126                          if (sendStatistics != null) {
127                              DataItem data = new DataItem("deletebench");
128                              data.setData(sendStatistics);
129                              mainService.addWork(data);
130                              success = true;
131                          }
132                      }
133                      if (command.equals("bulk")) {
134                          logger.info("'-c bulk' command supplied at startup");
135                          mainService.addWork(new DataItem("bulk"));
136                          success = true;
137                      }
138                      if (command.equals("collector")) {
139                          logger.info("'-c collector' command supplied at startup");
140                          String amount = (String) parameters.get(parameters
141                                  .indexOf("-a") + 1);
142                          if (amount != null) {
143                              DataItem data = new DataItem("collector");
144                              data.setCursorSize(Integer.parseInt(amount));
145                              mainService.addWork(data);
146                              success = true;
147                          }
148                      }
149                      if (command.equals("fillup")) {
150                          logger.info("'-c fillup' command supplied at startup");
151                          mainService.addWork(new DataItem("fillup"));
152                          success = true;
153                      }
154                  }
155              }
156
157          if (success == true) {
158              if (!parameters.isEmpty() && parameters.contains("-b")) {
159                  // this means we should run in the background with no GUI
160                  // starting
161                  // up
162                  logger.info("'-b' command supplied at startup."
163                          + " Run in background, no GUI starting up");
164                  success = true;
165              } else {
166                  Frame frame = new Frame();
167                  SplashScreen splashScreen = new SplashScreen("logo.jpg", frame);
168              }
169          } else {
170              logger.error("Wrong or no command parameters supplied at startup");
171              System.out.println();
172              System.out.println("ert usage:");
173              System.out.println();
174              System.out.println("    -c [join|first|user]   - command:");
175              System.out.println("        join"
176                      + "             join an existing network");
177              System.out.println("        first"
178                      + "            be the first node in a new network");
179              System.out.println("        user"
180                      + "             startup in user mode only,"
181                      + " no storage on this node");
182              System.out.println("    -b" + "                        "
183                      + "- runs the server in background (no GUI)");
184              System.out.println("    -v" + "                        "
185                      + "- prints out the application version");
186              System.out.println();
187              System.out.println("example: java -jar ert.jar -c join -b");
188              System.out.println();
189              System.exit(1);
190          }
191      }
192
193      /*
194       * (non-Javadoc)
```

```
195        *
196        * @see ch.nix.ert.ServiceInterface#registerAtDispatcher()
197        */
198      public void registerAtDispatcher() {
199          // TODO Auto-generated method stub
200          registerService("quit");
201      }
202
203      /*
204       * (non-Javadoc)
205       *
206       * @see java.lang.Runnable#run()
207       */
208      public void run() {
209          // TODO Auto-generated method stub
210          Thread.currentThread().setName("ERT.MainService");
211          while (true) {
212              DataItem dataItem = null;
213              dataItem = getWork();
214              // logger.debug("Got Work: "+dataItem);
215              String operation = dataItem.getOperationType();
216              if (operation.equals("quit")) {
217                  logger.info("quit request received, shutting down ERT");
218                  System.exit(0);
219              }
220          }
221      }
222
223 }
```

## File 37: main.SplashScreen

```
 1 /*
 2  * Created on Oct 20, 2004 5:12:51 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-26 15:28:00 +0200 (Tue, 26 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:     $Rev: 138 $
 8  * Location:    $URL: SplashScreen.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul written at
11  * the University of Zurich, Switzerland Department of Information Technology
12  * (www.ifi.unizh.ch) For copyright information please contact kussmaul@nix.ch
13  *
14  */
15
16 package ch.nix.ert.main;
17
18 import java.awt.BorderLayout;
19 import java.awt.Color;
20 import java.awt.Dimension;
21 import java.awt.Font;
22 import java.awt.Frame;
23 import java.awt.Toolkit;
24
25 import javax.swing.ImageIcon;
26 import javax.swing.JLabel;
27 import javax.swing.JWindow;
28 import javax.swing.SwingConstants;
29
30 /**
31  * Example code taken from: <br>
32  * http://www.javaworld.com/javaworld/javatips/jw-javatip104.html
33  *
34  * @author michi
35  */
36 class SplashScreen extends JWindow {
37
38      public SplashScreen(String filename, Frame f) {
39          super(f);
40          setBackground(Color.WHITE);
41
42          JLabel pictureLabel = new JLabel(new ImageIcon(SplashScreen.class.getResource(filename)));
43          JLabel versionLabel = new JLabel(" Version: $AntReplaceVersion$");
44          versionLabel.setFont(new Font("SansSerif",Font.PLAIN, 10));
45          versionLabel.setHorizontalAlignment(SwingConstants.LEFT);
46          getContentPane().add(pictureLabel, BorderLayout.CENTER);
47          getContentPane().add(versionLabel, BorderLayout.NORTH);
48          pack();
49          Dimension screenSize = Toolkit.getDefaultToolkit().getScreenSize();
50          Dimension labelSize = pictureLabel.getPreferredSize();
```

```
51              setLocation(screenSize.width / 2 - (labelSize.width / 2),
52                      screenSize.height / 2 - (labelSize.height / 2));
53          setVisible(true);
54          screenSize = null;
55          labelSize = null;
56      }
57 }
```

---

**File 38:** message.DataItem

```
 1 /*
 2  * Created on Oct 28, 2004 10:28:43 AM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-11-09 01:55:31 +0100 (Tue, 09 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 165 $
 8  * Location:    $URL: DataItem.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message;
17
18 import java.io.Serializable;
19 import java.util.ArrayList;
20
21 import ch.nix.ert.message.operations.OperationFactory;
22 import ch.nix.ert.message.operations.OperationInterface;
23
24 /**
25  * TODO Edit this text for a better description
26  *
27  * @author michi
28  */
29 public class DataItem implements Cloneable, Serializable {
30
31     protected Object[] operations = new Object[25];
32
33     public Object[] getInternalArray() {
34         return operations;
35     }
36
37     private void putValue(int opCode, Object value) {
38 //          if (operations[opCode] == null) {
39 //              operations[opCode] = OperationFactory.getOperation(opCode);
40 //          }
41         operations[opCode] = value;
42     }
43
44     private Object getValue(int opCode) {
45 //          if (operations[opCode] == null) {
46 //              return null;
47 //          }
48         return operations[opCode];
49     }
50
51     private void removeValue(int opCode) {
52         operations[opCode] = null;
53     }
54
55     private synchronized boolean existValue(int opCode) {
56         if (operations[opCode] == null) {
57             return false;
58         } else {
59             return true;
60         }
61     }
62
63     public Object clone() {
64         try {
65             Object dataItemCloned = super.clone();
66             ((DataItem) dataItemCloned).operations = (Object[]) operations.clone();
67             return dataItemCloned;
68         } catch (CloneNotSupportedException e) {
69             return null;
70         }
71     }
72
```

```
73      public DataItem getClone() {
74          return (DataItem) clone();
75      }
76
77      public DataItem() {
78      }
79
80      // helper Constructors, for common wrapping of Objects into a DataItem
81      public DataItem(String operation) {
82          this.setOperation(operation);
83      }
84
85      //    public DataItem(KeyPartition keyPartition) {
86      //        //TODO always make a deep copy of this object first, otherwise the
87      //        // keypartition does not allow concurrent access from multiple
88      //        // threads....
89      //        this.setKeyPartition(keyPartition);
90      //    }
91      //
92      //    public DataItem(String operation, KeyPartition keyPartition) {
93      //        //TODO always make a deep copy of this object first, otherwise the
94      //        // keypartition does not allow concurrent access from multiple
95      //        // threads....
96      //        this.setOperation(operation);
97      //        this.setKeyPartition(keyPartition);
98      //    }
99
100     public DataItem(InternalKey internalKey) {
101         this.setInternalKey(internalKey);
102     }
103
104     public DataItem(DataItem nestedItem) {
105         this.setNested(nestedItem);
106     }
107
108     public DataItem(String operation, DataItem nestedItem) {
109         this.setOperation(operation);
110         this.setNested(nestedItem);
111     }
112
113     public DataItem(Node destinationNode) {
114         this.setDirectDestination(destinationNode);
115     }
116
117     //for the most common keys accessors are supplied
118
119     public UserKey getUserKey() {
120         return (UserKey) this.getValue(OperationInterface.USERKEY);
121     }
122
123     public void setUserKey(UserKey key) {
124         this.putValue(OperationInterface.USERKEY, key);
125     }
126
127     public void removeUserKey() {
128         this.removeValue(OperationInterface.USERKEY);
129     }
130
131     // Destination InternalKey (where the DataItem needs to be processed)
132     public InternalKey getInternalKey() {
133         InternalKey internalKey = (InternalKey) this
134                 .getValue(OperationInterface.INTERNALKEY);
135         if (internalKey == null) {
136             UserKey userkey = this.getUserKey();
137             if (userkey == null) {
138                 return null;
139             }
140             internalKey = userkey.getInternalKey();
141         }
142         return internalKey;
143     }
144
145     public void setInternalKey(InternalKey iKey) {
146         this.putValue(OperationInterface.INTERNALKEY, iKey);
147     }
148
149     public void removeInternalKey() {
150         this.removeValue(OperationInterface.INTERNALKEY);
151     }
152
153     public String getNamespace() {
154         String namespace = (String) this.getValue(OperationInterface.NAMESPACE);
155         if (namespace == null) {
156             namespace = "default";
157         }
```

```
158        return namespace;
159    }
160
161    public void setNamespace(String namespace) {
162        this.putValue(OperationInterface.NAMESPACE, namespace);
163    }
164
165    public void removeNamespace() {
166        this.removeValue(OperationInterface.NAMESPACE);
167    }
168
169    public String getData() {
170        return (String) this.getValue(OperationInterface.DATA);
171    }
172
173    public void setData(String data) {
174        this.putValue(OperationInterface.DATA, data);
175    }
176
177    public void removeData() {
178        this.removeValue(OperationInterface.DATA);
179    }
180
181    public ArrayList getDataItems() {
182        return (ArrayList) this.getValue(OperationInterface.DATAITEMS);
183    }
184
185    public void addDataItems(DataItem data) {
186        if (this.getValue(OperationInterface.DATAITEMS) == null) {
187            //lazy initialisation of ArrayList
188            this.putValue(OperationInterface.DATAITEMS, new ArrayList());
189        }
190        ((ArrayList) this.getValue(OperationInterface.DATAITEMS)).add(data);
191    }
192
193    public void removeDataItems() {
194        this.removeValue(OperationInterface.DATAITEMS);
195    }
196
197    public String getTransactionId() {
198        return (String) this.getValue(OperationInterface.TRANSACTIONID);
199    }
200
201    public void setTransactionId(String transactionId) {
202        if (transactionId != null) {
203            this.putValue(OperationInterface.TRANSACTIONID, transactionId);
204        }
205    }
206
207    public void removeTransactionId() {
208        this.removeValue(OperationInterface.TRANSACTIONID);
209    }
210
211    public String getTransactionOperation() {
212        return (String) this.getValue(OperationInterface.TRANSACTIONOPERATION);
213    }
214
215    public void setTransactionOperation(String operation) {
216        if (operation != null) {
217            this.putValue(OperationInterface.TRANSACTIONOPERATION, operation);
218        }
219    }
220
221    public void removeTransactionOperation() {
222        this.removeValue(OperationInterface.TRANSACTIONOPERATION);
223    }
224
225    public Node getTransactionNodeId() {
226        return (Node) this.getValue(OperationInterface.TRANSACTIONNODEID);
227    }
228
229    public void setTransactionNodeId(Node node) {
230        if (node != null) {
231            this.putValue(OperationInterface.TRANSACTIONNODEID, node);
232        }
233    }
234
235    public void removeTransactionNodeId() {
236        this.removeValue(OperationInterface.TRANSACTIONNODEID);
237    }
238
239    public String getTransactionServiceId() {
240        return (String) this.getValue(OperationInterface.TRANSACTIONSERVICEID);
241    }
242
```

```
243      public void setTransactionServiceId(String serviceId) {
244          if (serviceId != null) {
245              this.putValue(OperationInterface.TRANSACTIONSERVICEID, serviceId);
246          }
247      }
248
249      public void removeTransactionServiceId() {
250          this.removeValue(OperationInterface.TRANSACTIONSERVICEID);
251      }
252
253      public String getTransactionType() {
254          String type = (String) this
255                  .getValue(OperationInterface.TRANSACTIONTYPE);
256          return type;
257      }
258
259      public void setTransactionType(String type) {
260          this.putValue(OperationInterface.TRANSACTIONTYPE, type);
261      }
262
263      public void removeTransactionType() {
264          this.removeValue(OperationInterface.TRANSACTIONTYPE);
265      }
266
267      public void setTransactionSplit(int splitSize) {
268          if (splitSize == 0) {
269              removeTransactionSplit();
270          } else {
271              this.putValue(OperationInterface.TRANSACTIONSPLIT, new Integer(
272                      splitSize));
273          }
274      }
275
276      public int getTransactionSplit() {
277          Integer value = (Integer) this
278                  .getValue(OperationInterface.TRANSACTIONSPLIT);
279          if (value == null) {
280              return 0;
281          } else {
282              return value.intValue();
283          }
284      }
285
286      public void removeTransactionSplit() {
287          this.removeValue(OperationInterface.TRANSACTIONSPLIT);
288      }
289
290      public boolean isTransactionLocal() {
291          Node node = getTransactionNodeId();
292          if (node == null) {
293              return true;
294          } else {
295              return getTransactionNodeId().isLocalhost();
296          }
297      }
298
299      public String getRouting() {
300          //values are: local, routing, direct
301          String routing = (String) this.getValue(OperationInterface.ROUTING);
302          return routing;
303      }
304
305      public void setRouting(String routingHint) {
306          this.putValue(OperationInterface.ROUTING, routingHint);
307      }
308
309      public void removeRounting() {
310          this.removeValue(OperationInterface.ROUTING);
311      }
312
313      public String getOperationType() {
314          String type = this.getTransactionType();
315          String operation = this.getOperation();
316          if (type != null && type.equals("end")) {
317              return operation + "-" + type;
318          }
319          return operation;
320      }
321
322      public String getOperation() {
323          return ((String) this.getValue(OperationInterface.OPERATION));
324      }
325
326      public void setOperation(String operation) {
327          this.putValue(OperationInterface.OPERATION, operation);
```

```
328        }
329
330        public void removeOperation() {
331            this.removeValue(OperationInterface.OPERATION);
332        }
333
334        public Node getDirectDestination() {
335            return (Node) this.getValue(OperationInterface.DIRECTDESTINATION);
336        }
337
338        public void setDirectDestination(Node node) {
339            this.putValue(OperationInterface.DIRECTDESTINATION, node);
340        }
341
342        public void removeDirectDestination() {
343            this.removeValue(OperationInterface.DIRECTDESTINATION);
344        }
345
346        public Node getDirectOriginator() {
347            return (Node) this.getValue(OperationInterface.DIRECTORIGINATOR);
348        }
349
350        public void setDirectOriginator(Node node) {
351            this.putValue(OperationInterface.DIRECTORIGINATOR, node);
352        }
353
354        public void removeDirectOriginator() {
355            this.removeValue(OperationInterface.DIRECTORIGINATOR);
356        }
357
358        public void setServiceOriginator(String service) {
359            this.putValue(OperationInterface.SERVICEORIGINATOR, service);
360        }
361
362        public String getServiceOriginator() {
363            return (String) this.getValue(OperationInterface.SERVICEORIGINATOR);
364        }
365
366        public void removeServiceOriginator() {
367            this.removeValue(OperationInterface.SERVICEORIGINATOR);
368        }
369
370        public DataItem getNested() {
371            return (DataItem) this.getValue(OperationInterface.NESTED);
372        }
373
374        public void setNested(DataItem data) {
375            this.putValue(OperationInterface.NESTED, data);
376        }
377
378        public void removeNested() {
379            this.removeValue(OperationInterface.NESTED);
380        }
381
382    //      // TODO currently only used for NodeDirectory to publish, unpublish
383    // data...
384    //      // will later perhaps removed to simplify the complex DataItem Object...
385    //      public KeyPartition getKeyPartition() {
386    //          return (KeyPartition) this.getValue("keypartition");
387    //      }
388    //
389    //      public void setKeyPartition(KeyPartition keyPartition) {
390    //          this.putValue("keypartition", keyPartition);
391    //      }
392    //
393    //      public void removeKeyPartition() {
394    //          this.removeValue("keypartition");
395    //      }
396
397        //START: Statistic related stuff
398        public void setDbSize(int size) {
399            this.putValue(OperationInterface.DBSIZE, new Integer(size));
400        }
401
402        public int getDbSize() {
403            return ((Integer) this.getValue(OperationInterface.DBSIZE)).intValue();
404        }
405
406        public void removeDbSize() {
407            this.removeValue(OperationInterface.DBSIZE);
408        }
409
410        public void setNodeLoad(double size) {
411            this.putValue(OperationInterface.NODELOAD, new Double(size));
412        }
```

```
413
414     public double getNodeLoad() {
415         return ((Double) this.getValue(OperationInterface.NODELOAD))
416                 .doubleValue();
417     }
418
419     public void removeNodeLoad() {
420         this.removeValue(OperationInterface.NODELOAD);
421     }
422
423     //END: Statistic related stuff
424
425     //START: Callback Service related stuff:
426     public void setDelay(long delay) {
427         this.putValue(OperationInterface.DELAY, new Long(delay));
428     }
429
430     public long getDelay() {
431         return ((Long) this.getValue(OperationInterface.DELAY)).longValue();
432     }
433
434     public void removeDelay() {
435         this.removeValue(OperationInterface.DELAY);
436     }
437
438     //END: Callback related stuff
439
440     public void setCursorSize(int size) {
441         this.putValue(OperationInterface.CURSORSIZE, new Integer(size));
442     }
443
444     public int getCursorSize() {
445         if (this.existValue(OperationInterface.CURSORSIZE)) {
446             return ((Integer) this.getValue(OperationInterface.CURSORSIZE))
447                     .intValue();
448         } else {
449             return 1;
450         }
451     }
452
453     public void removeCursorSize() {
454         this.removeValue(OperationInterface.CURSORSIZE);
455     }
456
457     public InternalKey getLastInternalKey() {
458         return (InternalKey) this.getValue(OperationInterface.LASTINTERNALKEY);
459     }
460
461     public void setLastInternalKey(InternalKey iKey) {
462         this.putValue(OperationInterface.LASTINTERNALKEY, iKey);
463     }
464
465     public void removeLastInternalKey() {
466         this.removeValue(OperationInterface.LASTINTERNALKEY);
467     }
468
469     public InternalKey getFirstInternalKey() {
470         return (InternalKey) this.getValue(OperationInterface.FIRSTINTERNALKEY);
471     }
472
473     public void setFirstInternalKey(InternalKey iKey) {
474         this.putValue(OperationInterface.FIRSTINTERNALKEY, iKey);
475     }
476
477     public void removeFirstInternalKey() {
478         this.removeValue(OperationInterface.FIRSTINTERNALKEY);
479     }
480
481     public void setKeyPartitionId(int id) {
482         this.putValue(OperationInterface.KEYPARTITIONID, new Integer(id));
483     }
484
485     public int getKeyPartitionId() {
486         if (this.existValue(OperationInterface.KEYPARTITIONID)) {
487             return ((Integer) this.getValue(OperationInterface.KEYPARTITIONID))
488                     .intValue();
489         } else {
490             return -1;
491         }
492     }
493
494     public void removeKeyPartitionId() {
495         this.removeValue(OperationInterface.KEYPARTITIONID);
496     }
497
```

```
498     public void setResult(boolean outcome) {
499         this.putValue(OperationInterface.RESULT, Boolean.valueOf(outcome));
500     }
501
502     public boolean getResult() {
503         Boolean value = (Boolean) this.getValue(OperationInterface.RESULT);
504         if (value == null) {
505             return true;
506         } else {
507             return value.booleanValue();
508         }
509     }
510
511     public void removeResult() {
512         this.removeValue(OperationInterface.RESULT);
513     }
514
515     public String toString() {
516         String output = "";
517         for (int i = 0; i < operations.length; i++) {
518             Object array_element = operations[i];
519             if (array_element != null) {
520                 output += OperationInterface.OPERATIONNAMES[i] + ":"
521                         + array_element + " ";
522             }
523         }
524         ;
525         return "[ " + output + "]";
526     }
527 }
```

---

**File 39:** message.DataItemCoder

---

```
 1 /*
 2  * Created on Oct 26, 2004 4:14:09 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-11-12 12:46:25 +0100 (Fri, 12 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 173 $
 8  * Location:    $URL: DataItemCoder.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message;
17
18 import java.nio.ByteBuffer;
19 import java.util.Enumeration;
20 import java.util.Hashtable;
21 import java.util.Iterator;
22
23 import org.apache.log4j.Logger;
24
25 import ch.nix.ert.communication.netty2.NodeCommunicationClientListener;
26 import ch.nix.ert.message.operations.OperationFactory;
27 import ch.nix.ert.message.operations.OperationInterface;
28
29 /**
30  * Codes the DataItem in ASN.1-kind way (key-length-value)
31  *
32  * @author michi
33  */
34 public class DataItemCoder {
35
36     //    // logger:
37     //    static Logger logger = Logger.getLogger(DataItemCoder.class);
38
39     private Object[] operations;
40
41     private int currentArrayElement = 0;
42
43     private byte[] bytePart = null;
44
45     private byte operationCode;
46
47     private boolean operationCodeRead = false;
48
49     private int offsetPart;
```

```
50
51      private boolean startNewPart = true;
52
53      final static private byte endMarker = (byte) 0xff;
54
55      private DataItem receivedDataItem = null;
56
57      public DataItemCoder() {
58          receivedDataItem = new DataItem();
59          operations = receivedDataItem.getInternalArray();
60      }
61
62      public void setDataItem(DataItem data) {
63          operations = data.getInternalArray();
64      }
65
66      public DataItem getDataItem() {
67          return receivedDataItem;
68      }
69
70      public synchronized boolean encode(ByteBuffer buf) {
71          int bufferSize = buf.remaining();
72          // first write all the data which could net yet be put on the buffer
73          if (bytePart != null) {
74              int stillToSend = bytePart.length - offsetPart;
75              if (bufferSize >= stillToSend) {
76                  buf.put(bytePart, offsetPart, stillToSend);
77                  bufferSize -= stillToSend;
78                  offsetPart = 0;
79              } else {
80                  buf.put(bytePart, offsetPart, bufferSize);
81                  offsetPart += bufferSize;
82                  return false;
83              }
84          }
85          while (true) {
86              for (int opCode = currentArrayElement; opCode < operations.length; opCode++) {
87                  Object element = operations[opCode];
88                  if (element != null) {
89                      OperationInterface operationCoder = OperationFactory
90                          .getOperation(opCode);
91                      bytePart = operationCoder.encode(element);
92                      int length = bytePart.length;
93                      if (bufferSize > 4) {
94                          buf.put((byte) opCode).putInt(length);
95                          bufferSize -= 5;
96                      }
97                      currentArrayElement = opCode + 1;
98                      if (bufferSize >= length) {
99                          buf.put(bytePart, 0, length);
100                         bufferSize -= length;
101                     } else {
102                         buf.put(bytePart, 0, bufferSize);
103                         offsetPart = bufferSize;
104                         return false;
105                     }
106                 }
107             }
108             if (bufferSize > 0) {
109                 buf.put(endMarker);
110                 //              byte[] buffer = buf.array();
111                 //              for (int i = 0; i < buffer.length; i++) {
112                 //                 System.out.println("buffer: "+buffer[i]);
113                 //              }
114                 // logger.debug("Wrote endmarker " + operations[6].getValue());
115                 return true;
116             } else {
117                 return false;
118             }
119         }
120     }
121
122     public synchronized boolean decode(ByteBuffer buf) {
123         int bufferSize = buf.remaining();
124         //System.out.println("decoder: "+buf);
125         while (true) {
126             // for (int opCode = currentArrayElement; opCode <
127             // operations.length; opCode++) {
128             if (startNewPart) {
129                 if (bufferSize >= 1 && !operationCodeRead) {
130                     operationCode = buf.get();
131                     operationCodeRead = true;
132                     bufferSize -= 1;
133                     if (operationCode == (byte) 0xff) {
134                         // this is the marker, that we have reached the end of
```

```
135                        // this dataItem
136                        return true;
137                    }
138                }
139                if (bufferSize > 4) {
140                    int sizeOfPart = buf.getInt();
141                    bytePart = new byte[sizeOfPart];
142                    offsetPart = 0;
143                    startNewPart = false;
144                    bufferSize -= 4;
145                } else {
146                    return false;
147                }
148            }
149            int stillToReceive = bytePart.length - offsetPart;
150            if (bufferSize >= stillToReceive) {
151                buf.get(bytePart, offsetPart, stillToReceive);
152                OperationInterface newOperation = OperationFactory
153                    .getOperation(operationCode);
154                operations[operationCode] = newOperation.decode(bytePart);
155                bufferSize -= stillToReceive;
156                operationCodeRead = false;
157                startNewPart = true;
158            } else {
159                buf.get(bytePart, offsetPart, bufferSize);
160                offsetPart += bufferSize;
161                return false;
162            }
163
164        }
165    }
166
167 }
```

---

### File 40: message.InternalKey

```
 1 /*
 2  * Created on Jul 27, 2004 4:14:37 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-11-01 11:26:55 +0100 (Mon, 01 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 151 $
 8  * Location:    $URL: InternalKey.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16
17 package ch.nix.ert.message;
18
19 import java.io.Serializable;
20
21 /**
22  *
23  * TODO Edit this text for a better description
24  * <p>
25  * TODO Has to check what the ranges of the keys are (e.g. 0-1000 or
26  * 0-100000000000) and should reject invalid values!
27  *
28  *
29  * @author michi
30  */
31 public class InternalKey implements Serializable, Comparable {
32
33     private long id;
34
35     static private long lowerLimit = 0;
36
37     static private long higherLimit = 1000000;
38
39     /**
40      *
41      */
42     public InternalKey(long uniqueId) {
43         uniqueId = adjustKey(uniqueId);
44         this.id = uniqueId;
45     }
46
```

```java
47      public InternalKey(String uniqueId) {
48          this.id = adjustKey(Long.decode(uniqueId).longValue());
49      }
50
51      public InternalKey getSubsequentInternalKey() {
52          return new InternalKey(this.id + 1);
53      }
54
55      public InternalKey getNthSubsequentInternalKey(long amount) {
56          return new InternalKey(this.id + amount);
57      }
58
59      public InternalKey getPreviousInternalKey() {
60          return new InternalKey(this.id - 1);
61      }
62
63      public InternalKey getNthPreviousInternalKey(long amount) {
64          return new InternalKey(this.id - amount);
65      }
66
67      /**
68       * @return Returns the uniqueID.
69       */
70      public long getId() {
71          return id;
72      }
73
74      static public long getHigherLimit() {
75          return higherLimit;
76      }
77
78      static public long getLowerLimit() {
79          return lowerLimit;
80      }
81
82      static public InternalKey getRandom() {
83          long difference = getHigherLimit() - getLowerLimit();
84          //Random generator = new Random();
85          long random = (long) (Math.random() * (double) difference);
86          return new InternalKey(random + getLowerLimit());
87      }
88
89      public UserKey getFirstUserKey() {
90          long internalKeyDiff = getHigherLimit() - getLowerLimit();
91          long userKeyDiff = UserKey.getHigherLimit()
92                  - UserKey.getLowerLimit();
93          float scalingFactor = (float) userKeyDiff / (float) internalKeyDiff;
94          //System.out.println("scalingFactor: "+scalingFactor);
95          long calculatedUserKey = (long) ((float) getId() * scalingFactor);
96          //System.out.println("calculatedUserKey: "+calculatedUserKey);
97          return new UserKey(calculatedUserKey
98                  + UserKey.getLowerLimit());
99      }
100
101     public UserKey getLastUserKey() {
102         long internalKeyDiff = getHigherLimit() - getLowerLimit();
103         long userKeyDiff = UserKey.getHigherLimit()
104                 - UserKey.getLowerLimit();
105         float scalingFactor = (float) userKeyDiff / (float) internalKeyDiff;
106         //System.out.println("scalingFactor: "+scalingFactor);
107         long calculatedUserKey = (long) ((float) (getId()+1) * scalingFactor);
108         //System.out.println("calculatedUserKey: "+calculatedUserKey);
109         return new UserKey(calculatedUserKey
110                 + UserKey.getLowerLimit() -1 );
111
112     }
113
114     /**
115      * @param uniqueID
116      *            The uniqueID to set.
117      */
118     public void setId(long uniqueId) {
119         uniqueId = adjustKey(uniqueId);
120         this.id = uniqueId;
121     }
122
123     /*
124      * (non-Javadoc)
125      *
126      * @see java.lang.Comparable#compareTo(java.lang.Object)
127      *
128      * Note: This class has a natural ordering that is inconsistent with equals
129      */
130     public int compareTo(Object o) {
131         InternalKey internalKey = (InternalKey) o;
```

```
132          if (internalKey.getId() > this.getId()) {
133              return -1;
134          } else if (internalKey.getId() < this.getId()) {
135              return 1;
136          } else {
137              return 0;
138          }
139      }
140
141      public long amountBetween(InternalKey internalKey) {
142          return Math.abs(this.getId() - internalKey.getId());
143      }
144
145      private long adjustKey(long uniqueId) {
146          if (uniqueId > higherLimit) {
147              return higherLimit;
148          }
149          if (uniqueId < lowerLimit) {
150              return lowerLimit;
151          }
152          return uniqueId;
153      }
154
155      public String toString() {
156          return String.valueOf(this.getId());
157      }
158 }
```

## File 41: message.Node

```
 1 /*
 2  * Created on Aug 2, 2004 11:06:57 AM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-28 22:38:52 +0200 (Thu, 28 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 144 $
 8  * Location:    $URL: Node.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16
17 package ch.nix.ert.message;
18
19 import java.io.Serializable;
20 import java.net.InetAddress;
21 import java.net.InetSocketAddress;
22 import java.net.UnknownHostException;
23
24 /**
25  * TODO Edit this text for a better description
26  *
27  * Represents the nodes (in this implementation IP/Port)
28  *
29  * @author michi
30  */
31 public class Node implements Serializable {
32
33      InetSocketAddress nodeAddress;
34
35      /**
36       *
37       */
38      public Node(String ip, int port) {
39          nodeAddress = new InetSocketAddress(ip, port);
40      }
41
42      public Node(InetAddress ip, int port) {
43          nodeAddress = new InetSocketAddress(ip, port);
44      }
45
46      /**
47       * Constructs a node with localhost as ip, and a  port as supplied
48       */
49      public Node(int port) {
50          try {
51              nodeAddress = new InetSocketAddress(InetAddress.getLocalHost(),port);
52          } catch (UnknownHostException e) {
```

```
 53                  // TODO Auto-generated catch block
 54                  e.printStackTrace();
 55          }
 56      }
 57
 58      /**
 59       * Constructs a node with localhost as ip, and a standard port as port... (15000)
 60       */
 61      public Node() {
 62          try {
 63              nodeAddress = new InetSocketAddress(InetAddress.getLocalHost(),15000);
 64          } catch (UnknownHostException e) {
 65              // TODO Auto-generated catch block
 66              e.printStackTrace();
 67          }
 68      }
 69
 70      /**
 71       * @return true: Node is on localhost
 72       *          <p>
 73       *          false: Node is on a remote host
 74       */
 75      public boolean isLocalhost() {
 76          InetAddress local;
 77          try {
 78              local = InetAddress.getLocalHost();
 79              return local.equals(nodeAddress.getAddress());
 80          } catch (UnknownHostException e) {
 81              // TODO Auto-generated catch block
 82              System.out.println("Error in Node isLocalhost: " + e);
 83          }
 84          return false;
 85      }
 86
 87 //    public boolean isAvailable() {
 88 //        // compose a "ping" packet!
 89 //        DataItem ping = new DataItem();
 90 //        ping.setDirectRecipient(this);
 91 //        ping.setOperation("ping");
 92 //        DataItem response = ServiceManager.getInstance().getMessageHandler().sendAndWait(ping);
 93 //        if (response == null) {
 94 //            return false;
 95 //        } else {
 96 //            return true;
 97 //        }
 98 //
 99 //    }
100
101      public InetSocketAddress getInetSocketAddress() {
102          return nodeAddress;
103      }
104
105      public String getIp() {
106          return nodeAddress.getAddress().getHostAddress();
107      }
108
109      public int getPort() {
110          return nodeAddress.getPort();
111      }
112
113      public String toString() {
114          return nodeAddress.getAddress().getHostAddress() + ":"
115                  + nodeAddress.getPort();
116      }
117
118 }
```

---

### File 42: message.UserKey

```
 1 /*
 2  * Created on Sep 14, 2004 12:40:18 PM
 3  *
 4  * Project:       ExtendendRootTree
 5  * Last Change: $Date: 2004-11-01 11:26:55 +0100 (Mon, 01 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 151 $
 8  * Location:    $URL: UserKey.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
```

```
14  *
15  */
16 package ch.nix.ert.message;
17
18 import java.io.Serializable;
19 import java.util.Comparator;
20
21 /**
22  * TODO Edit this text for a better description
23  *
24  * @author michi
25  */
26 public class UserKey implements Serializable, Comparable, Comparator {
27
28     private long id;
29
30     static private long lowerLimit = 0;
31
32     static private long higherLimit = 100000000;
33
34     public UserKey(long uniqueId) {
35         this.id = uniqueId;
36     }
37
38     public UserKey(Long uniqueId) {
39         this.id = uniqueId.longValue();
40     }
41
42     public UserKey(String uniqueId) {
43         this.id = Long.decode(uniqueId).longValue();
44     }
45
46     public UserKey getSubsequentUserKey() {
47         return new UserKey(this.id + 1);
48     }
49
50     public UserKey getPreviousInternalKey() {
51         return new UserKey(this.id - 1);
52     }
53
54     static public long getHigherLimit() {
55         return higherLimit;
56     }
57
58     static public long getLowerLimit() {
59         return lowerLimit;
60     }
61
62     public InternalKey getInternalKey() {
63         long internalKeyDiff = InternalKey.getHigherLimit()
64                 - InternalKey.getLowerLimit();
65         long userKeyDiff = getHigherLimit() - getLowerLimit();
66         float scalingFactor = (float) userKeyDiff / (float) internalKeyDiff;
67         long calculatedInternalKey = (long) ((float) getId() / scalingFactor);
68         return new InternalKey(calculatedInternalKey
69                 + InternalKey.getLowerLimit());
70     }
71
72     /**
73      * @return Returns the uniqueID.
74      */
75     public long getId() {
76         return id;
77     }
78
79     /*
80      * (non-Javadoc)
81      *
82      * @see java.lang.Comparable#compareTo(java.lang.Object)
83      */
84     public int compareTo(Object o) {
85         InternalKey internalKey = (InternalKey) o;
86         if (internalKey.getId() > this.getId()) {
87             return -1;
88         } else if (internalKey.getId() < this.getId()) {
89             return 1;
90         } else {
91             return 0;
92         }
93     }
94
95     public String toString() {
96         return String.valueOf(this.getId());
97     }
98
```

```
 99     public boolean equals(Object o) {
100         try {
101             UserKey userKey = (UserKey) o;
102             if (userKey.getId() == getId()) {
103                 return true;
104             } else {
105                 return false;
106             }
107         } catch (ClassCastException e) {
108             return false;
109         }
110     }
111
112     public int hashCode() {
113         throw new UnsupportedOperationException();
114     }
115
116     /*
117      * (non-Javadoc)
118      *
119      * @see java.util.Comparator#compare(java.lang.Object, java.lang.Object)
120      */
121     public int compare(Object o1, Object o2) {
122         UserKey userKey1 = (UserKey) o1;
123         UserKey userKey2 = (UserKey) o2;
124         return userKey1.compareTo(userKey2);
125     }
126
127 }
```

### File 43: message.operations.CursorSize

```
 1 /*
 2  * Created on Oct 28, 2004 3:25:48 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: CursorSize.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import netbyte.NetByte;
19 import netbyte.NetByteDecodingException;
20
21 /**
22  * TODO Edit this text for a better description
23  *
24  * @author michi
25  */
26 public class CursorSize implements OperationInterface {
27
28     public byte[] encode(Object cursorSize) {
29         return NetByte.encodeInt(((Integer) cursorSize).intValue());
30     }
31
32     public Object decode(byte[] buffer) {
33         try {
34             return new Integer(NetByte.decodeInt(buffer));
35         } catch (NetByteDecodingException e) {
36             // TODO Auto-generated catch block
37             e.printStackTrace();
38             return null;
39         }
40     }
41
42 }
```

### File 44: message.operations.Data

```
 1 /*
 2  * Created on Oct 28, 2004 3:18:41 PM
```

```
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: Data.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 /**
19  * TODO Edit this text for a better description
20  *
21  * @author michi
22  */
23 public class Data implements OperationInterface {
24
25     public byte[] encode(Object data) {
26         return ((String) data).getBytes();
27     }
28
29     public Object decode(byte[] buffer) {
30         return new String(buffer);
31     }
32
33 }
```

---

**File 45:** message.operations.DataItems

```
 1 /*
 2  * Created on Oct 28, 2004 3:19:08 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: DataItems.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 /**
19  * TODO Edit this text for a better description
20  *
21  * @author michi
22  */
23 public class DataItems implements OperationInterface {
24
25     /*
26      * (non-Javadoc)
27      *
28      * @see ch.nix.ert.message.operations.OperationInterface#encode()
29      */
30     public byte[] encode(Object DataItems) {
31         // TODO Auto-generated method stub
32         return null;
33     }
34
35     /*
36      * (non-Javadoc)
37      *
38      * @see ch.nix.ert.message.operations.OperationInterface#decode(byte[])
39      */
40     public Object decode(byte[] binary) {
41         // TODO Auto-generated method stub
42         return null;
43     }
44
45 }
```

**File 46:** message.operations.DbSize

```
 1 /*
 2  * Created on Oct 28, 2004 3:24:49 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: DbSize.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import netbyte.NetByte;
19 import netbyte.NetByteDecodingException;
20
21 /**
22  * TODO Edit this text for a better description
23  *
24  * @author michi
25  */
26 public class DbSize implements OperationInterface {
27
28     public byte[] encode(Object dbSize) {
29         return NetByte.encodeInt(((Integer) dbSize).intValue());
30     }
31
32     public Object decode(byte[] buffer) {
33         try {
34             return new Integer(NetByte.decodeInt(buffer));
35         } catch (NetByteDecodingException e) {
36             // TODO Auto-generated catch block
37             e.printStackTrace();
38             return null;
39         }
40     }
41 }
```

**File 47:** message.operations.Delay

```
 1 /*
 2  * Created on Oct 28, 2004 3:25:30 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: Delay.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import netbyte.NetByte;
19 import netbyte.NetByteDecodingException;
20
21 /**
22  * TODO Edit this text for a better description
23  *
24  * @author michi
25  */
26 public class Delay implements OperationInterface {
27
28     public byte[] encode(Object delay) {
29         return NetByte.encodeLong(((Long) delay).intValue());
30     }
31
32     public Object decode(byte[] buffer) {
33         try {
34             return new Long(NetByte.decodeLong(buffer));
35         } catch (NetByteDecodingException e) {
```

```
36              // TODO Auto-generated catch block
37              e.printStackTrace();
38              return null;
39          }
40      }
41 }
```

## File 48: message.operations.DirectDestination

```
 1 /*
 2  * Created on Oct 28, 2004 3:23:04 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:     $Rev: 150 $
 8  * Location:     $URL: DirectDestination.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import java.io.Serializable;
19 import java.net.InetAddress;
20 import java.net.InetSocketAddress;
21 import java.net.UnknownHostException;
22
23 import netbyte.NetByte;
24 import netbyte.NetByteDecodingException;
25 import ch.nix.ert.message.Node;
26
27 /**
28  * TODO Edit this text for a better description
29  *
30  * @author michi
31  */
32 public class DirectDestination implements OperationInterface {
33
34      public byte[] encode(Object value) {
35          Node destinationNode = (Node) value;
36          // will store it as:
37          //   ip
38          //   port
39          byte[] port = NetByte.encodeInt(destinationNode.getPort());
40          byte[] address = destinationNode.getInetSocketAddress().getAddress()
41              .getAddress();
42          byte[] returnValue = new byte[address.length + 4];
43          System.arraycopy(address, 0, returnValue, 0, address.length);
44          System.arraycopy(port, 0, returnValue, 4, 4);
45          //        for (int i = 0; i < returnValue.length; i++) {
46          //            System.out.println("ENCODE " + i + ": " + returnValue[i]);
47          //        }
48          return returnValue;
49      }
50
51      public Object decode(byte[] buffer) {
52          //        for (int i = 0; i < buffer.length; i++) {
53          //            System.out.println("DECODE " + i + ": " + buffer[i]);
54          //        }
55          try {
56              byte[] byteIp = new byte[4];
57              System.arraycopy(buffer, 0, byteIp, 0, 4);
58              InetAddress ip = InetAddress.getByAddress(byteIp);
59              int port = NetByte.decodeInt(buffer, 4);
60              //String ip = new String(buffer, 4, buffer.length - 4);
61              return new Node(ip, port);
62          } catch (Exception e) {
63              // TODO Auto-generated catch block
64              System.out.println("ERROR in here");
65              e.printStackTrace();
66              return null;
67          }
68      }
69 }
```

**File 49:** message.operations.DirectOriginator

```
 1 /*
 2  * Created on Oct 28, 2004 3:23:30 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-11-11 18:55:48 +0100 (Thu, 11 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 169 $
 8  * Location:    $URL: DirectOriginator.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import java.io.Serializable;
19 import java.net.InetAddress;
20
21 import netbyte.NetByte;
22 import ch.nix.ert.message.Node;
23
24 /**
25  * TODO Edit this text for a better description
26  *
27  * @author michi
28  */
29 public class DirectOriginator implements OperationInterface {
30
31     public byte[] encode(Object value) {
32         Node originatingNode = (Node) value;
33         // will store it as:
34         //   ip
35         //   port
36         byte[] port = NetByte.encodeInt(originatingNode.getPort());
37         byte[] address = originatingNode.getInetSocketAddress().getAddress()
38                 .getAddress();
39         byte[] returnValue = new byte[address.length + 4];
40         System.arraycopy(address, 0, returnValue, 0, address.length);
41         System.arraycopy(port, 0, returnValue, 4, 4);
42         return returnValue;
43     }
44
45     public Object decode(byte[] buffer) {
46         try {
47             byte[] byteIp = new byte[4];
48             System.arraycopy(buffer, 0, byteIp, 0, 4);
49             InetAddress ip = InetAddress.getByAddress(byteIp);
50             int port = NetByte.decodeInt(buffer, 4);
51             //String ip = new String(buffer, 4, buffer.length - 4);
52             return new Node(ip, port);
53         } catch (Exception e) {
54             // TODO Auto-generated catch block
55             System.out.println("ERROR in here");
56             e.printStackTrace();
57             return null;
58         }
59     }
60 }
```

**File 50:** message.operations.FirstInternalKey

```
 1 /*
 2  * Created on Oct 28, 2004 3:26:26 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: FirstInternalKey.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
```

```
17
18 import netbyte.NetByte;
19 import netbyte.NetByteDecodingException;
20
21 /**
22  * TODO Edit this text for a better description
23  *
24  * @author michi
25  */
26 public class FirstInternalKey implements OperationInterface {
27
28     public byte[] encode(Object iKey) {
29         return NetByte.encodeLong(((ch.nix.ert.message.InternalKey) iKey)
30                 .getId());
31     }
32
33     public Object decode(byte[] buffer) {
34         try {
35             return new ch.nix.ert.message.InternalKey(NetByte
36                     .decodeLong(buffer));
37         } catch (NetByteDecodingException e) {
38             // TODO Auto-generated catch block
39             e.printStackTrace();
40             return null;
41         }
42     }
43 }
```

---

**File 51:** message.operations.InternalKey

```
 1 /*
 2  * Created on Oct 28, 2004 3:18:18 PM
 3  *
 4  * Project:    ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: InternalKey.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import netbyte.NetByte;
19 import netbyte.NetByteDecodingException;
20
21 /**
22  * TODO Edit this text for a better description
23  *
24  * @author michi
25  */
26 public class InternalKey implements OperationInterface {
27
28     public byte[] encode(Object iKey) {
29         return NetByte.encodeLong(((ch.nix.ert.message.InternalKey) iKey).getId());
30     }
31
32     public Object decode(byte[] buffer) {
33         try {
34             return new ch.nix.ert.message.InternalKey(NetByte.decodeLong(buffer));
35         } catch (NetByteDecodingException e) {
36             // TODO Auto-generated catch block
37             e.printStackTrace();
38             return null;
39         }
40     }
41 }
```

---

**File 52:** message.operations.KeyPartitionId

```
 1 /*
 2  * Created on Oct 28, 2004 3:26:58 PM
 3  *
 4  * Project:    ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
```

```
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: KeyPartitionId.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import netbyte.NetByte;
19 import netbyte.NetByteDecodingException;
20
21 /**
22  * TODO Edit this text for a better description
23  *
24  * @author michi
25  */
26 public class KeyPartitionId implements OperationInterface {
27
28     public byte[] encode(Object partId) {
29         return NetByte.encodeInt(((Integer) partId).intValue());
30     }
31
32     public Object decode(byte[] buffer) {
33         try {
34             return new Integer(NetByte.decodeInt(buffer));
35         } catch (NetByteDecodingException e) {
36             // TODO Auto-generated catch block
37             e.printStackTrace();
38             return null;
39         }
40     }
41 }
```

---

### File 53: message.operations.LastInternalKey

```
 1 /*
 2  * Created on Oct 28, 2004 3:26:08 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: LastInternalKey.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import netbyte.NetByte;
19 import netbyte.NetByteDecodingException;
20
21 /**
22  * TODO Edit this text for a better description
23  *
24  * @author michi
25  */
26 public class LastInternalKey implements OperationInterface {
27
28     public byte[] encode(Object iKey) {
29         return NetByte.encodeLong(((ch.nix.ert.message.InternalKey) iKey)
30                 .getId());
31     }
32
33     public Object decode(byte[] buffer) {
34         try {
35             return new ch.nix.ert.message.InternalKey(NetByte
36                     .decodeLong(buffer));
37         } catch (NetByteDecodingException e) {
38             // TODO Auto-generated catch block
39             e.printStackTrace();
40             return null;
41         }
42     }
```

```
43 }
```

---

**File 54:** message.operations.Namespace

```
 1 /*
 2  * Created on Oct 28, 2004 2:58:08 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: Namespace.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 /**
19  * TODO Edit this text for a better description
20  *
21  * @author michi
22  */
23 public class Namespace implements OperationInterface {
24
25     public byte[] encode(Object namespace) {
26         return ((String) namespace).getBytes();
27     }
28
29     public Object decode(byte[] buffer) {
30         return new String(buffer);
31     }
32
33 }
```

---

**File 55:** message.operations.Nested

```
 1 /*
 2  * Created on Oct 28, 2004 3:24:26 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: Nested.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 /**
19  * TODO Edit this text for a better description
20  *
21  * @author michi
22  */
23 public class Nested implements OperationInterface {
24
25     public byte[] encode(Object nested) {
26         //TODO not supported yet...
27         return null;
28     }
29
30     public Object decode(byte[] buffer) {
31         return null;
32     }
33 }
```

---

**File 56:** message.operations.NodeLoad

```
 1 /*
 2  * Created on Oct 28, 2004 3:25:10 PM
```

```
 3   *
 4   * Project:      ExtendendRootTree
 5   * Last Change: $Date: 2004-11-09 01:55:31 +0100 (Tue, 09 Nov 2004) $
 6   * Changed by:   $Author: michi $
 7   * Revision:     $Rev: 165 $
 8   * Location:     $URL: NodeLoad.java $
 9   *
10   * This source-code is part of the diploma thesis of Michael Kussmaul
11   * written at the University of Zurich, Switzerland
12   * Department of Information Technology (www.ifi.unizh.ch)
13   * For copyright information please contact kussmaul@nix.ch
14   *
15   */
16  package ch.nix.ert.message.operations;
17
18  import netbyte.NetByte;
19  import netbyte.NetByteDecodingException;
20
21  /**
22   * TODO Edit this text for a better description
23   *
24   * @author michi
25   */
26  public class NodeLoad implements OperationInterface {
27
28      public byte[] encode(Object load) {
29          return NetByte.encodeDouble(((Double) load).doubleValue());
30      }
31
32      public Object decode(byte[] buffer) {
33          try {
34              return new Double(NetByte.decodeDouble(buffer));
35          } catch (NetByteDecodingException e) {
36              // TODO Auto-generated catch block
37              e.printStackTrace();
38              return null;
39          }
40      }
41  }
```

### File 57: message.operations.Operation

```
 1  /*
 2   * Created on Oct 26, 2004 5:18:41 PM
 3   *
 4   * Project:      ExtendendRootTree
 5   * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6   * Changed by:   $Author: michi $
 7   * Revision:     $Rev: 150 $
 8   * Location:     $URL: Operation.java $
 9   *
10   * This source-code is part of the diploma thesis of Michael Kussmaul
11   * written at the University of Zurich, Switzerland
12   * Department of Information Technology (www.ifi.unizh.ch)
13   * For copyright information please contact kussmaul@nix.ch
14   *
15   */
16  package ch.nix.ert.message.operations;
17
18  /**
19   * TODO Edit this text for a better description
20   *
21   * @author michi
22   */
23  public class Operation implements OperationInterface {
24
25      public byte[] encode(Object operationCode) {
26          return ((String) operationCode).getBytes();
27      }
28
29      public Object decode(byte[] buffer) {
30          return  new String(buffer);
31      }
32
33      //    public boolean encode(Object operation, ByteBuffer buf) {
34      //        String code = (String) operation;
35      //        byte[] raw = ((String) operation).getBytes();
36      //        int length = raw.length;
37      //        if (buf.remaining() < length) {
38      //            return false;
39      //        }
40      //        buf.put(OPERATION).putInt(length).put(raw);
```

```
41    //           return true;
42    //      }
43    //
44    //    public void decode (DataItem data, ByteBuffer buf) {
45    //          int length = buf.getInt();
46    //          byte[] raw = new byte[length];
47    //          String opCode = new String(buf.get(raw).array());
48    //          data.putValue("operation", opCode);
49    //      }
50
51      /*
52       * (non-Javadoc)
53       *
54       * @see ch.nix.ert.message.operations.OperationInterface#getName()
55       */
56
57 }
```

## File 58: message.operations.OperationFactory

```
 1 /*
 2  * Created on Oct 28, 2004 12:21:26 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-10-28 17:44:43 +0200 (Thu, 28 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 143 $
 8  * Location:    $URL: OperationFactory.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import org.apache.log4j.Logger;
19
20 import ch.nix.ert.keypartition.MigrationClient;
21
22 import com.sleepycat.je.tree.DIN;
23
24 /**
25  * TODO Edit this text for a better description
26  *
27  * @author michi
28  */
29 public class OperationFactory {
30
31     // logger:
32     static Logger logger = Logger.getLogger(OperationFactory.class);
33
34     public static OperationInterface getOperation(int opCode) {
35         switch (opCode) {
36         case OperationInterface.OPERATION:
37             return new Operation();
38         case OperationInterface.USERKEY:
39             return new UserKey();
40         case OperationInterface.INTERNALKEY:
41             return new InternalKey();
42         case OperationInterface.NAMESPACE:
43             return new Namespace();
44         case OperationInterface.DATA:
45             return new Data();
46         case OperationInterface.DATAITEMS:
47             return new DataItems();
48         case OperationInterface.TRANSACTIONID:
49             return new TransactionId();
50         case OperationInterface.TRANSACTIONOPERATION:
51             return new TransactionOperation();
52         case OperationInterface.TRANSACTIONNODEID:
53             return new TransactionNodeId();
54         case OperationInterface.TRANSACTIONSERVICEID:
55             return new TransactionServiceId();
56         case OperationInterface.TRANSACTIONTYPE:
57             return new TransactionType();
58         case OperationInterface.TRANSACTIONSPLIT:
59             return new TransactionSplit();
60         case OperationInterface.ROUTING:
61             return new Routing();
62         case OperationInterface.DIRECTDESTINATION:
```

```
63              return new DirectDestination();
64          case OperationInterface.DIRECTORIGINATOR:
65              return new DirectOriginator();
66          case OperationInterface.SERVICEORIGINATOR:
67              return new ServiceOriginator();
68          case OperationInterface.NESTED:
69              return new Nested();
70          case OperationInterface.DBSIZE:
71              return new DbSize();
72          case OperationInterface.NODELOAD:
73              return new NodeLoad();
74          case OperationInterface.DELAY:
75              return new Delay();
76          case OperationInterface.CURSORSIZE:
77              return new CursorSize();
78          case OperationInterface.LASTINTERNALKEY:
79              return new LastInternalKey();
80          case OperationInterface.FIRSTINTERNALKEY:
81              return new FirstInternalKey();
82          case OperationInterface.KEYPARTITIONID:
83              return new KeyPartitionId();
84          case OperationInterface.RESULT:
85              return new Result();
86          default:
87              logger.error("Can't handle unknown Message Attribute with ID"+opCode);
88              return null;
89          }
90      }
91
92 }
```

### File 59: message.operations.OperationInterface

```
 1 /*
 2  * Created on Oct 26, 2004 5:22:53 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-11-02 01:36:34 +0100 (Tue, 02 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 154 $
 8  * Location:    $URL: OperationInterface.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import java.io.Serializable;
19
20 /**
21  * TODO Edit this text for a better description
22  *
23  * @author michi
24  */
25 public interface OperationInterface extends Serializable {
26
27      // list of all keys for DataItem
28      public static final int OPERATION = 0;
29
30      public static final int USERKEY = 1;
31
32      public static final int INTERNALKEY = 2;
33
34      public static final int NAMESPACE = 3;
35
36      public static final int DATA = 4;
37
38      public static final int DATAITEMS = 5;
39
40      public static final int TRANSACTIONID = 6;
41
42      public static final int TRANSACTIONOPERATION = 7;
43
44      public static final int TRANSACTIONNODEID = 8;
45
46      public static final int TRANSACTIONSERVICEID = 9;
47
48      public static final int TRANSACTIONTYPE = 10;
49
```

```
50      public static final int TRANSACTIONSPLIT = 11;
51
52      public static final int ROUTING = 12;
53
54      public static final int DIRECTDESTINATION = 13;
55
56      public static final int DIRECTORIGINATOR = 14;
57
58      public static final int SERVICEORIGINATOR = 15;
59
60      public static final int NESTED = 16;
61
62      public static final int DBSIZE = 17;
63
64      public static final int NODELOAD = 18;
65
66      public static final int DELAY = 19;
67
68      public static final int CURSORSIZE = 20;
69
70      public static final int FIRSTINTERNALKEY = 21;
71
72      public static final int LASTINTERNALKEY = 22;
73
74      public static final int KEYPARTITIONID = 23;
75
76      public static final int RESULT = 24;
77
78      public static final String[] OPERATIONNAMES = { "Operation", "UserKey",
79              "InternalKey", "Namespace", "Data", "DataItems", "TransactionId",
80              "TransactionOperation", "TransactionNodeId",
81              "TransactionServiceId", "TransactionType", "TransactionSplit",
82              "Routing", "DirectDestination", "DirectOriginator",
83              "ServiceOriginator", "Nested", "DbSize", "NodeLoad", "Delay",
84              "CursorSize", "FirstInternalKey", "LastInternalKey", "KeyPartitionId", "Result" };
85
86      public byte[] encode(Object value);
87
88      public Object decode(byte[] binary);
89
90      // public boolean encode(Object value, ByteBuffer buf);
91
92      // public void decode(DataItem data, ByteBuffer buf);
93
94      //    public Object getValue();
95      //
96      //    public void setValue(Object value);
97      //
98      //    public String getName();
99
100 }
```

---

**File 60:** message.operations.Result

```
1  /*
2   * Created on Oct 28, 2004 3:27:17 PM
3   *
4   * Project:      ExtendendRootTree
5   * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
6   * Changed by:  $Author: michi $
7   * Revision:    $Rev: 150 $
8   * Location:    $URL: Result.java $
9   *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18
19 /**
20  * TODO Edit this text for a better description
21  *
22  * @author michi
23  */
24 public class Result implements OperationInterface {
25
26      public byte[] encode(Object result) {
27          return new byte[(byte) ( ((Boolean) result).booleanValue() ? 1 : 0)];
28      }
```

```
29
30     public Object decode(byte[] buffer) {
31         return Boolean.valueOf(buffer[0] == 0);
32     }
33 }
```

**File 61:** message.operations.Routing

```
1  /*
2   * Created on Oct 28, 2004 3:22:37 PM
3   *
4   * Project:     ExtendendRootTree
5   * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
6   * Changed by:  $Author: michi $
7   * Revision:    $Rev: 150 $
8   * Location:    $URL: Routing.java $
9   *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 /**
19  * Values for routing are: <br>
20  * local: menans this dataItem has a local destination <br>
21  * routing: means NodeDirectory needs to resolve destination <br>
22  * direct: means IP is already know, where to send this dataItem <br>
23  *
24  * @author michi
25  */
26 public class Routing implements OperationInterface {
27
28     public byte[] encode(Object routing) {
29         return ((String) routing).getBytes();
30     }
31
32     public Object decode(byte[] buffer) {
33         return new String(buffer);
34     }
35
36 }
```

**File 62:** message.operations.ServiceOriginator

```
1  /*
2   * Created on Oct 28, 2004 3:23:59 PM
3   *
4   * Project:     ExtendendRootTree
5   * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
6   * Changed by:  $Author: michi $
7   * Revision:    $Rev: 150 $
8   * Location:    $URL: ServiceOriginator.java $
9   *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 /**
19  * TODO Edit this text for a better description
20  *
21  * @author michi
22  */
23 public class ServiceOriginator implements OperationInterface {
24
25     public byte[] encode(Object serviceOriginator) {
26         return ((String) serviceOriginator).getBytes();
27     }
28
29     public Object decode(byte[] buffer) {
30         return new String(buffer);
31     }
32
         public Object decode(byte[] buffer) {
```

```
33 }
```

---

**File 63:** message.operations.TransactionId

```
1 /*
2  * Created on Oct 28, 2004 3:19:39 PM
3  *
4  * Project:      ExtendendRootTree
5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
6  * Changed by:  $Author: michi $
7  * Revision:    $Rev: 150 $
8  * Location:    $URL: TransactionId.java $
9  *
10 * This source-code is part of the diploma thesis of Michael Kussmaul
11 * written at the University of Zurich, Switzerland
12 * Department of Information Technology (www.ifi.unizh.ch)
13 * For copyright information please contact kussmaul@nix.ch
14 *
15 */
16 package ch.nix.ert.message.operations;
17
18 /**
19  * TODO Edit this text for a better description
20  *
21  * @author michi
22  */
23 public class TransactionId implements OperationInterface {
24
25     public byte[] encode(Object transactionId) {
26         return ((String) transactionId).getBytes();
27     }
28
29     public Object decode(byte[] buffer) {
30         return new String(buffer);
31     }
32 }
```

---

**File 64:** message.operations.TransactionNodeId

```
1 /*
2  * Created on Oct 28, 2004 3:20:25 PM
3  *
4  * Project:      ExtendendRootTree
5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
6  * Changed by:  $Author: michi $
7  * Revision:    $Rev: 150 $
8  * Location:    $URL: TransactionNodeId.java $
9  *
10 * This source-code is part of the diploma thesis of Michael Kussmaul
11 * written at the University of Zurich, Switzerland
12 * Department of Information Technology (www.ifi.unizh.ch)
13 * For copyright information please contact kussmaul@nix.ch
14 *
15 */
16 package ch.nix.ert.message.operations;
17
18 import java.net.InetAddress;
19
20 import netbyte.NetByte;
21 import ch.nix.ert.message.Node;
22
23 /**
24  * TODO Edit this text for a better description
25  *
26  * @author michi
27  */
28 public class TransactionNodeId implements OperationInterface {
29
30     public byte[] encode(Object value) {
31         Node transactionNodeId = (Node) value;
32         // will store it as:
33         //   ip
34         //   port
35         byte[] port = NetByte.encodeInt(transactionNodeId.getPort());
36         byte[] address = transactionNodeId.getInetSocketAddress().getAddress()
37                 .getAddress();
38         byte[] returnValue = new byte[address.length + 4];
39         System.arraycopy(address, 0, returnValue, 0, address.length);
40         System.arraycopy(port, 0, returnValue, 4, 4);
```

```
41          return returnValue;
42      }
43
44      public Object decode(byte[] buffer) {
45          try {
46              byte[] byteIp = new byte[4];
47              System.arraycopy(buffer, 0, byteIp, 0, 4);
48              InetAddress ip = InetAddress.getByAddress(byteIp);
49              int port = NetByte.decodeInt(buffer, 4);
50              //String ip = new String(buffer, 4, buffer.length - 4);
51              return new Node(ip, port);
52          } catch (Exception e) {
53              // TODO Auto-generated catch block
54              System.out.println("ERROR in here");
55              e.printStackTrace();
56              return null;
57          }
58      }
59
60 }
```

**File 65:** message.operations.TransactionOperation

```
 1 /*
 2  * Created on Oct 28, 2004 3:20:04 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: TransactionOperation.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 /**
19  * TODO Edit this text for a better description
20  *
21  * @author michi
22  */
23 public class TransactionOperation implements OperationInterface {
24
25      public byte[] encode(Object transactionOp) {
26          return ((String) transactionOp).getBytes();
27      }
28
29      public Object decode(byte[] buffer) {
30          return new String(buffer);
31      }
32
33 }
```

**File 66:** message.operations.TransactionServiceId

```
 1 /*
 2  * Created on Oct 28, 2004 3:20:49 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: TransactionServiceId.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 /**
19  * TODO Edit this text for a better description
20  *
```

```
21  * @author michi
22  */
23 public class TransactionServiceId implements OperationInterface {
24
25     public byte[] encode(Object transactionServiceId) {
26         return ((String) transactionServiceId).getBytes();
27     }
28
29     public Object decode(byte[] buffer) {
30         return new String(buffer);
31     }
32 }
```

---

### File 67: message.operations.TransactionSplit

```
 1 /*
 2  * Created on Oct 28, 2004 3:21:51 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: TransactionSplit.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import netbyte.NetByte;
19 import netbyte.NetByteDecodingException;
20
21 /**
22  * TODO Edit this text for a better description
23  *
24  * @author michi
25  */
26 public class TransactionSplit implements OperationInterface {
27
28     public byte[] encode(Object transactionSplit) {
29         return NetByte.encodeInt(((Integer) transactionSplit).intValue());
30     }
31
32     public Object decode(byte[] buffer) {
33         try {
34             return new Integer(NetByte.decodeInt(buffer));
35         } catch (NetByteDecodingException e) {
36             // TODO Auto-generated catch block
37             e.printStackTrace();
38             return null;
39         }
40     }
41
42 }
```

---

### File 68: message.operations.TransactionType

```
 1 /*
 2  * Created on Oct 28, 2004 3:21:15 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: TransactionType.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 /**
19  * TODO Edit this text for a better description
```

```
20  *
21  * @author michi
22  */
23 public class TransactionType implements OperationInterface {
24
25     public byte[] encode(Object transactionType) {
26         return ((String) transactionType).getBytes();
27     }
28
29     public Object decode(byte[] buffer) {
30         return new String(buffer);
31     }
32
33 }
```

**File 69:** message.operations.UserKey

```
 1 /*
 2  * Created on Oct 28, 2004 3:17:54 PM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-10-31 12:28:25 +0100 (Sun, 31 Oct 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 150 $
 8  * Location:    $URL: UserKey.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.message.operations;
17
18 import netbyte.NetByte;
19 import netbyte.NetByteDecodingException;
20
21 /**
22  * TODO Edit this text for a better description
23  *
24  * @author michi
25  */
26 public class UserKey implements OperationInterface {
27
28     public byte[] encode(Object uKey) {
29         return NetByte.encodeLong(((ch.nix.ert.message.UserKey) uKey).getId());
30     }
31
32     public Object decode(byte[] buffer) {
33         try {
34             return new ch.nix.ert.message.UserKey(NetByte.decodeLong(buffer));
35         } catch (NetByteDecodingException e) {
36             // TODO Auto-generated catch block
37             e.printStackTrace();
38             return null;
39         }
40     }
41
42 }
```

**File 70:** user.UserService

```
 1 /*
 2  * Created on Oct 5, 2004 11:18:43 AM
 3  *
 4  * Project:     ExtendendRootTree
 5  * Last Change: $Date: 2004-11-07 22:22:16 +0100 (Sun, 07 Nov 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 164 $
 8  * Location:    $URL: UserService.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.user;
17
```

```
18 import java.util.Date;
19 import java.util.StringTokenizer;
20
21 import org.apache.log4j.Logger;
22
23 import ch.nix.ert.base.ServiceBase;
24 import ch.nix.ert.message.DataItem;
25 import ch.nix.ert.message.InternalKey;
26 import ch.nix.ert.message.Node;
27 import ch.nix.ert.message.UserKey;
28
29 /**
30  * TODO Edit this text for a better description
31  *
32  * @author michi
33  */
34 public class UserService extends ServiceBase {
35
36     //  logger:
37     static Logger logger = Logger.getLogger(UserService.class);
38
39     static private UserService instance = null;
40
41     static int reqSec = 0;
42
43     static int throughput = 0;
44
45     static int collectorAmount = 0;
46
47     static int collected = 0;
48
49     static String[][] statistics;
50
51     private UserService() {
52     }
53
54     static public synchronized UserService getInstance() {
55         if (null == instance) {
56             instance = new UserService();
57             logger.info("UserService started");
58         }
59         return instance;
60     }
61
62     /*
63      * (non-Javadoc)
64      *
65      * @see ch.nix.ert.ServiceInterface#registerAtDispatcher()
66      */
67     public void registerAtDispatcher() {
68         // nothing to register, as this UserService should not listen to any
69         // keywords, he should only retrieve transaction results
70         registerService("fillup");
71         registerService("bulk");
72         registerService("writebench");
73         registerService("readbench");
74         registerService("deletebench");
75         registerService("collect");
76         registerService("collector");
77     }
78
79     /*
80      * (non-Javadoc)
81      *
82      * @see java.lang.Runnable#run()
83      */
84     public void run() {
85         Thread.currentThread().setName("ERT.UserService");
86         long startMillis = 0;
87         long startMillis2 = 0;
88         int amount = 2000;
89         int counter = amount;
90         boolean bulkMode = false;
91         while (true) {
92             DataItem dataItem = null;
93             dataItem = getWork();
94             // logger.debug("Got Work: "+dataItem);
95             String operation = dataItem.getOperationType();
96
97             if (operation.equals("collector")) {
98                 // I'm a collector, so I will retrieve the amount of machines, i
99                 // have to collect...
100                 collectorAmount = dataItem.getCursorSize();
101                 statistics = new String[collectorAmount][10];
102             } else if (operation.equals("collect")) {
```

```
103                     String statisticData = dataItem.getData();
104                     StringTokenizer st = new StringTokenizer(statisticData);
105                     int i = 0;
106                     while (st.hasMoreTokens()) {
107                         statistics[collected][i] = st.nextToken();
108                         i++;
109                     }
110                     collected++;
111                     if (collected >= collectorAmount) {
112                         // now generate statistics
113                         for (int j = 0; j < 10; j++) {
114                             String timeslice = "";
115                             double timesliceSum = 0;
116                             for (int k = 0; k < collectorAmount; k++) {
117                                 String number = statistics[k][j];
118                                 int len = number.length();
119                                 if (len < 4) {
120                                     for (int index = 0; index < (4 - len); index++) {
121                                         number = " " + number;
122                                     }
123                                 }
124                                 timeslice += number + " ";
125                                 timesliceSum += Double.parseDouble(number);
126                             }
127                             // timesliceSum = timesliceSum / (double) collectorAmount;
128                             String sum = "" + Math.round(timesliceSum);
129                             int len = sum.length();
130                             if (len < 9) {
131                                 for (int index = 0; index < (9 - len); index++) {
132                                     sum = " " + sum;
133                                 }
134                             }
135                             logger.error("Timeslice: " + sum + " - " + timeslice);
136                         }
137                         collected = 0;
138                     }
139                 } else if (operation.equals("bulk")) {
140                     bulkMode = true;
141                     try {
142                         Thread.sleep(10000);
143                     } catch (InterruptedException e) {
144                         // TODO Auto-generated catch block
145                         e.printStackTrace();
146                     }
147                     logger.info("Starting benchmark");
148                     Date start = new Date();
149                     startMillis = start.getTime();
150                     for (int i = 0; i < amount; i = i + 1) {
151                         long random = (long) (Math.random() * (double) 100000000);
152                         put(random, "Fillup " + random);
153                         //System.out.println(random);
154                     }
155                     Date stop = new Date();
156                     long stopMillis = stop.getTime();
157                     logger
158                             .info("Sending "
159                                     + amount
160                                     + " entries took: "
161                                     + (stopMillis - startMillis)
162                                     + "ms ("
163                                     + ((double) amount
164                                             / (double) (stopMillis - startMillis) * (double) 1000)
165                                     + " requests/sec)");
166
167                 } else if (operation.equals("put-end") && bulkMode) {
168                     // logger.info("Got: "+dataItem);
169                     if (counter == amount) {
170                         Date start2 = new Date();
171                         startMillis2 = start2.getTime();
172                     }
173                     counter-;
174                     //                    if (counter%1 == 0) {
175                     //                        logger.info("Still missing responses: "+counter);
176                     //                    }
177                     if (counter < 1) {
178                         Date stop = new Date();
179                         long stopMillis = stop.getTime();
180                         logger
181                                 .info("Receiving "
182                                         + amount
183                                         + " entries took: "
184                                         + (stopMillis - startMillis2)
185                                         + "ms ("
186                                         + ((double) amount
187                                                 / (double) (stopMillis - startMillis2) * (double) 1000)
```

```
188                                    + " requests/sec)");
189                    logger
190                            .info("Overall Storing "
191                                    + amount
192                                    + " entries took: "
193                                    + (stopMillis - startMillis)
194                                    + "ms ("
195                                    + ((double) amount
196                                            / (double) (stopMillis - startMillis) * (double) 1000)
197                                    + " requests/sec)");
198                    // restart
199                    counter = amount;
200                    addWork(new DataItem("bulk"));
201                }
202            }
203
204        else if (operation.equals("writebench")) {
205            final Node collectorNode = new Node(dataItem.getData(), 0);
206            final int dataSize = dataItem.getDbSize();
207            try {
208                Thread.sleep(10000);
209            } catch (InterruptedException e) {
210                // TODO Auto-generated catch block
211                e.printStackTrace();
212            }
213            Thread t1 = new Thread("ThreadT1") {
214                long startMillis;
215                public void run() {
216
217                    int round = 1;
218                    StringBuffer buffer = new StringBuffer(60);
219                    String data = "";
220                    for (int i = 0; i < dataSize; i++) {
221                        data += "d";
222                    }
223                    while (true) {
224                        Date start = new Date();
225                        startMillis = start.getTime();
226                        for (int i = 0; i < throughput + 10; i = i + 1) {
227                            long random = (long) (Math.random() * (double) 100000000);
228                            put(random, data);
229                        }
230                        Date stop = new Date();
231                        long stopMillis = stop.getTime();
232                        long timeToWait = 1000 - (stopMillis - startMillis) % 1000;
233                        try {
234                            Thread.sleep(timeToWait);
235                        } catch (InterruptedException e) {
236                            // TODO Auto-generated catch block
237                            e.printStackTrace();
238                        }
239                        logger.info("Received requests in last second: "
240                                + reqSec);
241                        buffer.append(" "+reqSec);
242                        if (round >= 10) {
243                            // send statistics to collector:
244                            DataItem collectorData = new DataItem("collect");
245                            collectorData.setDirectDestination(collectorNode);
246                            collectorData.setRouting("direct");
247                            collectorData.setData(buffer.toString());
248                            addWork(collectorData);
249                            round = 0;
250                            buffer = new StringBuffer(60);
251                        }
252                        round++;
253                        throughput = reqSec;
254                        reqSec = 0;
255                    }
256                }
257            };
258            logger.info("Starting benchmark");
259            t1.start();
260        } else if (operation.equals("put-end")) {
261            reqSec++;
262        } else if (operation.equals("fillup")) {
263            try {
264                Thread.sleep(10000);
265            } catch (InterruptedException e) {
266                // TODO Auto-generated catch block
267                e.printStackTrace();
268            }
269            Thread t1 = new Thread("ThreadT1") {
270                long startMillis;
271
272                public void run() {
```

```
273                              logger.info("Starting filling DB");
274                              for (int i = 0; i < 1000000; i = i + 1) {
275                                  put(i, "Fillup " + i);
276                                  if (i % 1000 == 0) {
277                                      logger.info("Filling DB: " + i
278                                              + " entries in db");
279                                  }
280                                  try {
281                                      Thread.sleep(10);
282                                  } catch (InterruptedException e) {
283                                      // TODO Auto-generated catch block
284                                      e.printStackTrace();
285                                  }
286                              }
287                              logger.info("Stopped filling DB: 1'000'000 entries");
288                          }
289                      };
290                      t1.start();
291                  } else if (operation.equals("deletebench")) {
292                      final Node collectorNode = new Node(dataItem.getData(), 0);
293                      try {
294                          Thread.sleep(10000);
295                      } catch (InterruptedException e) {
296                          // TODO Auto-generated catch block
297                          e.printStackTrace();
298                      }
299                      Thread t1 = new Thread("ThreadT1") {
300                          long startMillis;
301
302                          public void run() {
303                              int round = 1;
304                              StringBuffer buffer = new StringBuffer(60);
305                              while (true) {
306                                  Date start = new Date();
307                                  startMillis = start.getTime();
308                                  for (int i = 0; i < throughput + 10; i = i + 1) {
309                                      long random = (long) (Math.random() * (double) 100000000);
310                                      delete(random);
311                                  }
312                                  Date stop = new Date();
313                                  long stopMillis = stop.getTime();
314                                  long timeToWait = 1000 - (stopMillis - startMillis) % 1000;
315                                  try {
316                                      Thread.sleep(timeToWait);
317                                  } catch (InterruptedException e) {
318                                      // TODO Auto-generated catch block
319                                      e.printStackTrace();
320                                  }
321                                  logger.info("Received requests in last second: "
322                                          + reqSec);
323                                  buffer.append(" "+reqSec);
324                                  if (round >= 10) {
325                                      // send statistics to collector:
326                                      DataItem collectorData = new DataItem("collect");
327                                      collectorData.setDirectDestination(collectorNode);
328                                      collectorData.setRouting("direct");
329                                      collectorData.setData(buffer.toString());
330                                      addWork(collectorData);
331                                      round = 0;
332                                      buffer = new StringBuffer(60);
333                                  }
334                                  round++;
335                                  throughput = reqSec;
336                                  reqSec = 0;
337                              }
338                          }
339                      };
340                      logger.info("Starting benchmark");
341                      t1.start();
342                  } else if (operation.equals("delete-end")) {
343                      reqSec++;
344                  } else if (operation.equals("readbench")) {
345                      final Node collectorNode = new Node(dataItem.getData(), 0);
346                      try {
347                          Thread.sleep(10000);
348                      } catch (InterruptedException e) {
349                          // TODO Auto-generated catch block
350                          e.printStackTrace();
351                      }
352                      Thread t1 = new Thread("ThreadT1") {
353                          long startMillis;
354
355                          public void run() {
356                              int round = 1;
357                              StringBuffer buffer = new StringBuffer(60);
```

```
358                             while (true) {
359                                 Date start = new Date();
360                                 startMillis = start.getTime();
361                                 for (int i = 0; i < throughput + 10; i = i + 1) {
362                                     long random = (long) (Math.random() * (double) 100000000);
363                                     get(random);
364                                 }
365                                 Date stop = new Date();
366                                 long stopMillis = stop.getTime();
367                                 long timeToWait = 1000 - (stopMillis - startMillis) % 1000;
368                                 try {
369                                     Thread.sleep(timeToWait);
370                                 } catch (InterruptedException e) {
371                                     // TODO Auto-generated catch block
372                                     e.printStackTrace();
373                                 }
374                                 logger.info("Received requests in last second: "
375                                         + reqSec);
376                                 buffer.append(" "+reqSec);
377                                 if (round >= 10) {
378                                     // send statistics to collector:
379                                     DataItem collectorData = new DataItem("collect");
380                                     collectorData.setDirectDestination(collectorNode);
381                                     collectorData.setRouting("direct");
382                                     collectorData.setData(buffer.toString());
383                                     addWork(collectorData);
384                                     round = 0;
385                                     buffer = new StringBuffer(60);
386                                 }
387                                 round++;
388                                 throughput = reqSec;
389                                 reqSec = 0;
390                             }
391                         }
392                     };
393                     logger.info("Starting benchmark");
394                     t1.start();
395                 } else if (operation.equals("get-end")) {
396                     reqSec++;
397                 }

399         }

401     }

403     public void sendRaw(DataItem dataItem) {
404         transactionBegin(dataItem);
405     }

407     public synchronized void put(long key, String data) {
408         UserKey userKey = new UserKey(key);
409         InternalKey internalKey = userKey.getInternalKey();

411         DataItem putData = new DataItem();
412         putData.setUserKey(userKey);
413         putData.setInternalKey(internalKey);
414         putData.setData(data);
415         putData.setOperation("put");
416         putData.setRouting("routing");
417         transactionBegin(putData);
418     }

420     public synchronized boolean putSync(long key, String data) {
421         put(key, data);

423         // now wait for a response...
424         DataItem result = getWork();
425         return result.getResult();
426     }

428     public synchronized void get(long key) {
429         UserKey userKey = new UserKey(key);
430         InternalKey internalKey = userKey.getInternalKey();

432         DataItem getData = new DataItem();
433         getData.setUserKey(userKey);
434         getData.setInternalKey(internalKey);
435         getData.setOperation("get");
436         getData.setRouting("routing");
437         transactionBegin(getData);
438     }

440     public synchronized Object getSync(long key) {
441         get(key);
442
```

```
443          // now wait for a response...
444          DataItem result = getWork();
445          return result.getData();
446      }
447
448      public synchronized void delete(long key) {
449          UserKey userKey = new UserKey(key);
450          InternalKey internalKey = userKey.getInternalKey();
451
452          DataItem deleteData = new DataItem();
453          deleteData.setUserKey(userKey);
454          deleteData.setInternalKey(internalKey);
455
456          deleteData.setOperation("delete");
457          deleteData.setRouting("routing");
458          transactionBegin(deleteData);
459      }
460
461      public synchronized boolean deleteSync(long key) {
462          delete(key);
463
464          // now wait for a response...
465          DataItem result = getWork();
466          return result.getResult();
467      }
468
469 }
```

---

### File 71: util.LinkedHashMap

```
 1 /*
 2  * Created on Sep 24, 2004 5:33:57 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-09-24 17:58:24 +0200 (Fri, 24 Sep 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 97 $
 8  * Location:    $URL: LinkedHashMap.java $
 9  *
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.util;
17
18 import java.util.Map;
19
20 /**
21  * TODO Edit this text for a better description
22  *
23  * @author michi
24  */
25 public class LinkedHashMap extends java.util.LinkedHashMap {
26
27      public void setMaxSize(int size) {
28          MAX_ENTRIES = size;
29      }
30
31      private int MAX_ENTRIES = 10000;
32
33      protected boolean removeEldestEntry(Map.Entry eldest) {
34          return size() > MAX_ENTRIES;
35      }
36
37 }
```

---

### File 72: util.Queue

```
 1 /*
 2  * Created on Sep 8, 2004 5:08:21 PM
 3  *
 4  * Project:      ExtendendRootTree
 5  * Last Change: $Date: 2004-09-08 17:26:50 +0200 (Wed, 08 Sep 2004) $
 6  * Changed by:  $Author: michi $
 7  * Revision:    $Rev: 79 $
 8  * Location:    $URL: Queue.java $
 9  *
```

```
10  * This source-code is part of the diploma thesis of Michael Kussmaul
11  * written at the University of Zurich, Switzerland
12  * Department of Information Technology (www.ifi.unizh.ch)
13  * For copyright information please contact kussmaul@nix.ch
14  *
15  */
16 package ch.nix.ert.util;
17
18 import EDU.oswego.cs.dl.util.concurrent.LinkedQueue;
19
20 /**
21  * TODO Edit this text for a better description
22  *
23  * @author michi
24  */
25 public class Queue extends LinkedQueue {
26
27     private String myName = "noname";
28
29     private Object lockObject = new Object();
30
31     /**
32      *
33      */
34     public Queue() {
35         super();
36     }
37
38     public Queue(String name) {
39         super();
40         myName = name;
41     }
42
43     public void setName(String name) {
44         synchronized (lockObject) {
45             myName = name;
46         }
47     }
48
49     public synchronized String getName() {
50         synchronized (lockObject) {
51             return myName;
52         }
53     }
54
55 }
```