



Facoltà di Ingegneria

Corso di Laurea Specialistica in Ingegneria Informatica

Tesi di Laurea

In-Memory Database: Competitive Landscape and Performance Analysis

Laureando

Valerio Barbagallo

Relatore

Prof. Paolo Merialdo

Universit Roma Tre

Tutor Aziendale

Dr. Michele Aiello

Applications Director & Co. Founder

ERIS4 s.r.l.

Anno Accademico 2007-2008

Dedicato a te

Acknowledgements

Thank you for your attention!

Abstract

Naturally this is a joke. The abstract must be written after all chapters are written. Here i'm going to tell you what i'm talking about. The thesis is divided into the following chapters:

1. In-Memory Database

- A brief introduction
 - Definition
 - Usage
- IMDBs overview
 -
 -

2. Database Performance Analysis

-
-

3. High availability with in-memory database?

Contents

Acknowledgements	ii
Abstract	iii
I Competitive Landscape	1
1 A Brief Introduction	2
1.1 Definition	2
1.2 History	3
1.3 Application Scenario	3
1.4 Comparison against Traditional DBMS	4
1.4.1 Caching	5
1.4.2 Data-Transfer Overhead	6
1.4.3 Transaction Processing	6
1.5 ACID Support: Adding Durability	6
1.5.1 On-Line Backup	7
1.5.2 Transaction Logging	7
1.5.3 High Availability Implementation	8

1.5.4	NVRAM	9
2	In-Memory Database Overview	10
2.1	The Analysis' Structure	10
2.1.1	Common Advantages	11
2.1.2	Common Disadvantages	11
2.1.3	Common references	12
2.2	Prevayler	12
2.2.1	Advantages	13
2.2.2	Disadvantages	14
2.2.3	Project Info	15
2.2.4	Usage	15
2.3	HSQLDB	21
2.3.1	Advantages	21
2.3.2	Disadvantages	22
2.3.3	Project Info	22
2.4	Db4o	22
2.4.1	Advantages	23
2.4.2	Disadvantages	23
2.4.3	Project Info	23
2.5	Hxsq	23
2.5.1	Advantages	24
2.5.2	Disadvantages	24
2.5.3	Project Info	24
2.6	H2	24

2.6.1	Advantages	25
2.6.2	Disadvantages	25
2.6.3	Project Info	25
2.7	Derby	25
2.7.1	Advantages	26
2.7.2	Disadvantages	26
2.7.3	Project Info	26
2.8	SQLite	26
2.8.1	Advantages	27
2.8.2	Disadvantages	27
2.8.3	Project Info	27
2.9	Firebird	27
2.9.1	Advantages	28
2.9.2	Disadvantages	28
2.9.3	Project Info	28
2.10	MySQL	28
2.10.1	Advantages	29
2.10.2	Disadvantages	29
2.10.3	Project Info	29
2.11	ExtremeDB	29
2.12	Polyhedra	30
2.13	TimesTen	30
2.14	Csql	30
2.15	SolidDB	30
2.16	MonetDB	31

2.17 RDM Embedded	31
2.18 FastDB	31
2.19 QuiLogic	32
2.20 Pico4	32
2.21 Pico4v2	32
2.22 Conclusion	32
 II Performance Analysis	 33
 3 Performance Analysis Introduction	 34
3.1 Impartiality Problem	34
3.1.1 Benchmark History	35
3.1.2 Benchmark Differences	36
3.1.3 The Axiom	36
3.2 Measures	37
3.3 Choosing a Database	39
 4 Test Suite	 41
4.1 Why Define Test Scenarios	41
4.2 Base Test Case	42
4.2.1 Race Test: Timed	43
4.2.2 Race Test: Transactions	44
4.3 Load Test Case	45
4.3.1 Real Time Prepaid System	46
4.4 Acid Test Case	51

5	Database Benchmark Softwares Overview	53
5.1	Benchmark Requirements	53
5.2	The Open Source Database Benchmark	53
5.3	Transaction Processing Performance Council	53
5.4	Apache JMeter	54
5.5	Poleposition	54
6	The New In-Memory Database Benchmark Application	55
6.1	Functional View	55
6.2	Development View	55
6.3	Plug-In Architecture	55
7	Results' Analysis	56
A	The first appendix	57

List of Figures

List of Tables

Part I

Competitive Landscape

Chapter 1

A Brief Introduction

“While familiar on desktops and servers, databases are a recent arrival to embedded systems. Like any organism dropped into a new environment, databases must evolve. A new type of DBMS, the in-memory database system (IMDS), represents the latest step in DBMSes’ adaptation to embedded systems” [Graves 02].

In this chapter we are going to make a brief introduction to the in-memory databases, explaining what they are, their use, their strength and weakness.

1.1 Definition

An in-memory database (IMDB), also called main memory database (MMDB), in-memory database system (IMDS) or real-time database (RTDB), is a database management system that relies on main memory for data storage. While the bandwidth of hard disks is just 1 order of magnitude slower than the main memory’s bandwidth, the disk access time is about 3 order of

magnitude slower than the RAM access time, and thus in-memory databases can be much more faster than traditional database management systems (DBMS).

1.2 History

Initially embedded systems developers produced their own data management solutions. But with the market competition requiring smarter devices, applications with expanding feature set will have to manage increasingly complex data structures. As a consequence, these data management solutions were outgrowing, and became difficult to maintain and extend.

Therefore embedded systems developers turned to commercial databases. But the first embedded databases were not the ideal solution. They were traditional DBMS with complex caching logic to increase performance, and with a lot of unnecessary features for the device that make use of embedded databases. Furthermore these features cause the application to exceed available memory and CPU resources.

In-memory databases have emerged specifically to meet the performance needs and resource availability in embedded systems.

1.3 Application Scenario

Often in-memory databases run as embedded database, but it's not their only use. Thanks to their high performance, these databases are particularly useful for all that kind of applications that need fast access to the data. Some

examples:

- real time applications which don't need to be persisted either because it doesn't change, or the data can be reconstructed: imagine a routing table of a router with millions of record and data access in less than few milliseconds; the routing table can be rebuilt [Fowler 05].
- real time applications with durability needs which capture, analyze and respond intelligently to important events, requiring high performance in terms of throughput and mainly latency (traditional DBMS can be clustered to increase the throughput, but with no great benefits in terms of latency). Infact almost all IMDBs can be persistent on disk while still keeping higher performance compared to traditional DBMSs.
- in-memory databases are also very useful for developers of traditional database systems for testing purpose: in a enterprise application running a test suite can take long; switching to an IMDB can reduce the whole build time of the application.

1.4 Comparison against Traditional DBMS

In-memory databases eliminate disk I/O and exist only in RAM, but they are not simply a traditional database loaded into main memory. Linux systems already have the capability to create a RAM disk, a file system in main memory. But a traditional database deployed in a such virtual hard disk doesn't provide the same benefits of a pure IMDB. In-memory databases are

less complex than a traditional DBMS fully deployed in RAM, and lead to a minor usage of CPU and RAM.

Comparing IMDBs with traditional databases we can find at least 3 key differences:

- Caching.
- Data-transfer overhead.
- Transaction processing.

1.4.1 Caching

All traditional DBMS software incorporates caching mechanisms to keep the most used records in main memory to reduce the performance issue introduced by the disk latency. The removal of caching brings to the elimination of the following tasks:

- cache synchronization, used to keep the portion of the database loaded in main memory consistent with the physical database image.
- cache lookup, a task that handle the cached page, determining if the data requested is in cache

Therefore, removing the cache, IMDBs eliminate a great source of complexity and performance overhead, reducing the work for the CPU and, comparing to a traditional DBMS fully deployed in main memory, also the RAM.

1.4.2 Data-Transfer Overhead

Traditional DBMS adds a remarkable data transfer overhead due not only to the DB cache, but to the file system and his cache too. In contrast an IMDBs, which eliminate these steps, have little or no data transfer.

There is also no need for the application to make a copy of the data in local memory, because IMDBs give to the application a pointer to the data that reside in the database. In this way the data is accessed through the database API that protect the data itself.

1.4.3 Transaction Processing

In an hard-disk database the recovery process, from a disaster or a transaction abort, is based on a log file, that is update every time a transaction is executed. To provide transactional integrity, IMDBs maintain a before image of the objects updated and in case of a transaction abort, the before image are restored in a very efficient way. Therefore another complex, memory-intensive task is eliminated from the IMDB, unless the need to add durability to the system, because there is no reason to keep transaction log files.

1.5 ACID Support: Adding Durability

When we choose a database we expect from it to provide ACID support: atomicity, consistency, isolation and durability. While the first three features are usually supported by in-memory databases, pure IMDBs, in their simplest

form, don't provide durability: main memory is a volatile memory and loses all stored data during a reset.

With their high performance, IMDBs are a good solution for time-critical applications, but when the durability need arises they may not seem a proper solution anymore. To achieve durability [Gorine 04], in-memory database systems can use several solutions:

- On-Line Backup.
- Transaction logging.
- High availability implementations.
- Non volatile RAM (NVRAM).

1.5.1 On-Line Backup

On-line backup is a backup performed while the database is on-line and available for read/write. This is the simplest solution, but offer a minimum degree of durability.

1.5.2 Transaction Logging

A transaction log is a history of actions executed by a database management system. To guarantee ACID properties over crashes or hardware failures the log must be written in a non-volatile media, usually an hard disk. If the system fails and is restarted, the database image can be restored from this log file.

The recovery process acts similarly to traditional DBMS with a roll-back or a roll-forward recovery. Checkpoint (or snapshot) can be used to speed up this process. However this technique implies the usage of persistence memory such as an hard disk, that is a bottleneck, especially during the resume of the database.

Although this aspect, IMDBs are still faster than traditional DBMS:

1. transaction logging requires exactly one write to the file system, while disk-based databases not only need to write on the log, but also the data and the indexes (and even more writes with larger transactions).
2. transaction logging may be usually set to different level of transaction durability. A trade-off between performance and durability is allowed by setting the log mechanism to be synchronous or asynchronous.

1.5.3 High Availability Implementation

High availability is a system design protocol and associated implementation: in this case it is a database replication, with automatic fail over to an identical standby database. A replicated database consists of failure-independent nodes, making sure data is not lost even in case of a node failure. This is an effective way to achieve database transaction durability.

In a simile way to transaction logging, a trade-off between performance and durability is achieved by setting the replication as eager (synchronous) or lazy (asynchronous).

1.5.4 NVRAM

To achieve durability a IMDB can support non volatile ram (NVRAM): usually a static RAM backed up with battery power, or some sort of EEPROM. In this way the DBMS can recover the data also after a reboot.

This is a very attractive durability option for IMDBs. NVRAM in contrast to transaction logging and database replication does not involve disk I/O latency and neither the communication overhead.

Despite this, vendors rarely provide support for this technology. One of the major problems to this approach is the limited write-cycles of this kind of memory, such as a flash memory. On the other hand there is some new memory device that has been proposed to address this problem, but the cost of such devices is rather expensive, in particular considering the huge amount of memory needed by IMDBs.

Chapter 2

In-Memory Database Overview

There is a variety of in-memory database systems which can be used to maintain a database in main memory, both commercial and open source. Although all of them share the capability to maintain the database in main memory, they offer different sets of feature.

In this chapter we will make a competitive landscape of the most famous in-memory databases.

2.1 The Analysis' Structure

Every in-memory database will be analyzed investigating their advantages and disadvantages, the stability and reliability of the project and how much development is going on. Eventually we will go to the development stage and, in some case, we will also have a look "under the hood" to see how the DB works.

2.1.1 Common Advantages

Not all IMDBs have the same advantages, but generally they share some feature:

Lightweight is one of the common advantages IMDBs share. This kind of databases is really simple: they don't implement any of the complex mechanism used by traditional database to speed up the disk I/O. Thus, most of the time, an in-memory database consist of a simple jar, whose size is less than 1 MB.

Robustness is a direct consequence of the previous point: (citing Henry Ford) it is impossible to break something that doesn't exist.

High Performance is another common feature for every IMDB, because they all store the whole database in main memory, avoiding disk access.

2.1.2 Common Disadvantages

On the other side, all IMDBs require high quantity of main memory. It is larger as the database image increases, but this doesn't mean that every IMDB uses the same RAM quantity for the same database image.

A common question is about what may happen when the RAM is not enough. Although this question is really interesting, generally an IMDB makes the assumption there is always enough main memory to hold the database image.

2.1.3 Common references

The source of informations described in the following sections comes from the web, mainly from sourceforge or ohloh.net. Both these web site offer a variety of news such as the size of developer team, the frequency of commits, the date of the last realese, etc.

Another source is, naturally, the IMDBs' official web site. Although this is a huge source, it cannont be considered impartial. Therefore the following analysis may not be consistent with facts.

2.2 Prevayler

Prevayler is an object persistence library for Java, written in Java. It keeps data in main memory and any change is written to a journal file for system recovery. This is an implementation of an architectural style that is called System Prevalence by the Prevayler developer team.

Therefore Prevayler, being an object prevalence layer, provides transparent persistence for Plain Old Java Objects. In the prevalent model, the object data is kept in memory in native, language-specific object format, and therefore avoid any marshalling to an RDBMS or other data storage system. To avoid losing data and to provide durability a snapshot of data is regularly saved to disk and all changes are serialized to a log of transactions which is also stored on disk [Wuestefeld 01].

All is based on the Prevalent Hypothesis: that there is enough RAM to hold all business objects in your system.

2.2.1 Advantages

Prevayler is a lightweight java library, just 350 KB, that is extremely simple to use. There is no separate database server to run. With Prevayler you can program with real objects, there is no use of SQL, there is no impedance mismatch such as when programming with RDBMS. Moreover Prevayler doesn't require the domain objects to implement or extend any class in order to be persistent (except `java.io.Serializable`, but this is only a marker interface).

It is very fast. Simply keeping objects in memory in their language-specific format is both orders of magnitude faster and more programmer-friendly than the multiple conversions that are needed when the objects are stored and retrieved from an RDBMS. The only disk access is the streaming of the transactions to the journal file that should be about one thousand¹ transactions per second on an average desktop computer.

The thread safety is guaranteed. Actually, in the default Prevayler implementation all writes to the system are synchronized. One write at a time. So there's no threading issues at all. Therefore there is no more multithreading issue such as locking, atomicity, consistency and isolation.

Finally Prevayler supports the execution only in RAM, like a pure in-memory database should work. But it can also provide persistence through a journal file, as we described above. Moreover Prevayler supports snapshots of the database's image, which serves to speed up the database's boot, and server replication, enabling query load-balancing and system fault-tolerance [Prevayler 08]. This last feature is really promising, because in-

¹This is what Prevayler team says on their official web site on Mar 26, 2008.

memory databases suffer the start up process, allowing it to be used in an enterprise application. Although this feature is not ready yet: see paragraph 2.2.4 for further details at page 19.

2.2.2 Disadvantages

On the other hand of Prevayler's simplicity, there is some restriction due to the fact that only changes are written in the journal file: transaction must be completely deterministic and repeatable in order to recover the state of the object (for instance it's not possible to use `System.currentTimeMillis()`).

In addition, when Prevayler is embedded in your application, the only client of your application's data is the application itself [Hobbs 03]. While deploying Prevayler as a server, the access to the data is still limited to whom who speaks the host language (you can't use another programming language unless you make your own serialization mechanism) and, at the same time, knows the model objects. For all these reasons there are no administration and migration tools.

Another problem is related to the Prevalent Hypothesis. If, for any reason, the RAM is not enough and an object model is swapped out of main memory, Prevayler will become very slow, more than traditional RDBMS's [Miller 03]. In fact Prevayler doesn't use any mechanism to optimize the disk I/O, such as traditional DBMS's mechanisms (eg: indexing, caching etc.) Anyway this issue belongs to all pure in-memory databases.

2.2.3 Project Info

Klaus Wuestefeld is the founder and main developer of Prevayler, an open source project. This project started in 2001, from what sourceforge reports, and had a great development until 2005. The community is still active, but the last update is dated at 25/05/2007 when version 2.3 (the latest) was released. The development team is composed of 8 developers, including Klaus Wuestefeld.. The current development status is declared to be production/stable.

2.2.4 Usage

In this example, and in all the followings, we are going to use a simple Plain Old Java Object as business object that need to be persisted: `Number`. It has one single private field, which is the value of the number itself, and the relative getter and setter methods.

```
public class Number{
    private int value;
    public Number(){
    }
    public Number(int value) {
        this.value = value;
    }
    public int getValue() {
        return value;
    }
    public void setValue(int number) {
        this.value = number;
    }
}
```

```
}
```

Prevayler usage is quite different from a traditional RDBMS with JDBC driver. Prevayler is accessed through a native driver, and it acts similarly to a framework: your business objects must implement the interface `java.io.Serializable` in order to be persisted (which is only a marker interface); and all modifications to these objects must be encapsulated in transaction objects. Therefore our business object will appear as follow:

```
public class Number implements Serializable { ...
```

Basic: main memory and transaction logging

This example is about the default Prevayler's usage: how to create an in-memory database without losing durability. This property is achieved with a transaction log file, whose usage is totally transparent. To initialize Prevayler Database you need to create a Prevayler, providing it with the directory where you want to save your database, and the object which is going to be saved in this Prevayler database. This step can be compared to a `CREATE TABLE` in SQL, but only one object will be saved in your Prevayler database. Therefore, from a RDBMS perspective where the table is a class and each row is one instance, you may want to initialize Prevayler with a `List` or a `Map` of `Number`. Here is an example:

```
public class Main{  
    private static final String DIRECTORY_DB = "numberDB";  
    public static void main(String[] args) throws Exception {  
        Prevayler prevayler = PrevaylerFactory.createPrevayler(  
            new HashMap<Integer, Number>(), DIRECTORY_DB);
```

```

    new Main().fillPrevaylerDb(prevayler);
    new Main().readPrevaylerDb(prevayler);
}
...

```

It's important to understand that the state of the object you use to create this database will not be saved in the database itself. To insert any `Number` in the database a transaction must be executed:

```

private void fillPrevaylerDb(Prevayler prevayler) throws \
    InterruptedException {
    for (int i = 0; i<100; i++){
        prevayler.execute(new InsertNumberTransaction(new Number\
            (i)));
        System.out.println("The value of the number inserted is \
            =" + i);
    }
}

```

The parameter of the method `execute` must be a class that extends `org.prevayler.Transaction`. Every insert, update or delete (in other words: any write operation) requires to be executed inside a transaction. In this particular case this is the the code:

```

public class InsertNumberTransaction implements Transaction{
    private Number number;
    public InsertNumberTransaction(Number number) {
        this.number = number;
    }
    public void executeOn(Object prevalentSystem, Date ignored\
        ) {

```

```

    Map<Integer, Number> map = (Map<Integer, Number>) \
        prevalentSystem;
    map.put(number.getValue(), number);
}
}

```

It's important to note that when you stop the database, and then you restart it, Prevayler will execute all the transactions exactly the same number of times they were executed before shutting down the process ². While a snapshot should avoid this behavior.

Finally, reading from Prevayler database is really simple and doesn't require any transaction:

```

private void readPrevaylerDb(Prevayler prevayler) {
    Map<Integer, Number> map = (Map<Integer, Number>) prevayler.\
        prevalentSystem();
    Set<Integer> keys = map.keySet();
    for (Integer key : keys) {
        Number number = map.get(key);
        System.out.println("Reading the number " + number.\
            getValue());
    }
}
}

```

Only RAM

This example show how to make Prevayler run only in main memory, without the writes to the journal file, in the case you want a database even faster and you don't care for durability or you don't have write permission. There is

²This is the reason why the transactions must be deterministic

only one line of code which need to be modified to be able to run Prevayler in such a way:

```
Prevayler prevayler = PrevaylerFactory.\n    createTransientPrevayler(new HashMap<Integer, Number>());
```

Instead of creating a persistent prevayler, you just need to `createTransientPrevayler`. You can also decide to create the transient prevayler from a snapshot, and then work only in main memory: really useful for the execution of your test case. But you can't take a snapshot while using transient prevayler, therefore you need to disable the snapshot in your code when switching from the default to the transient prevayler, otherwise a `IOException` will raise.

With this method you have no durability, but it is even faster than the first example, about 10 times faster. And moreover it is more scalable, because there is no more bottleneck caused by the hard disk.

Server Replication

Prevayler can support also a server replication modality. Also in this case, only a small change to the database initialization is needed. Quite obviously you need two different kind of initialization: server side and clients side.

As regards the server, you need to specify the port number and, then, simply create the object Prevayler:

```
public class MainServer {
    private static final String DB_DIRECTORY_PATH = "\n    numberReplicaDB";
    private static final int PORT_NUMBER = 37127;
    public static void main(String[] args) throws Exception {
        PrevaylerFactory factory = new PrevaylerFactory();
```

```

    factory.configurePrevalentSystem(new HashMap<Integer,\
        Number>());
    factory.configurePrevalenceDirectory(DB_DIRECTORY_PATH);
    factory.configureReplicationServer(PORT_NUMBER);
    Prevayler prevayler = factory.create();
    ...
    // execute your transactions
    ...
    // The server will continue to listen/run for incoming \
        connections
    ...
}
}

```

As for the clients, you have to tell Prevayler not only the port number, but the ip address too. Only one line of code changes from the server:

```

public class MainReplicant {
    private static final String DB_DIRECTORY_NAME = "numberDB"\
        ;
    private static final int PORT_NUMBER = 37127;
    private static final String SERVER_IP = "10.0.2.2";

    public static void main(String[] args) throws Exception {
        PrevaylerFactory factory = new PrevaylerFactory();
        factory.configurePrevalentSystem(new HashMap<Integer,\
            Number>());
        factory.configurePrevalenceDirectory(DB_DIRECTORY_NAME);
        factory.configureReplicationClient(SERVER_IP, \
            PORT_NUMBER);
        Prevayler prevayler = factory.create();
    }
}

```


2.3.2 Disadvantages

2.3.3 Project Info

2.4 Db4o

22

advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages advantages
advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages

disadvantages disadvantages disadvantages disadvantages disadvantages dis-
advantages disadvantages disadvantages disadvantages disadvantages disad-
vantages disadvantages disadvantages disadvantages disadvantages disadvan-
tages

info
info
info
info info info info info info info info info info info info

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages advantages
advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages

disadvantages disadvantages disadvantages disadvantages disadvantages dis-
advantages disadvantages disadvantages disadvantages disadvantages disad-
vantages disadvantages disadvantages disadvantages disadvantages disadvan-
tages

info
info
info
info info info info info info info info info info info info

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages advantages
advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages

disadvantages disadvantages disadvantages disadvantages disadvantages dis-
advantages disadvantages disadvantages disadvantages disadvantages disad-
vantages disadvantages disadvantages disadvantages disadvantages disadvan-
tages

info info info info info info info info info info info info info info info info
info info info info info info info info info info info info info info info info
info info info info info info info info info info info info info info info info
info info info info info info info info info info info

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages advantages
advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages

disadvantages disadvantages disadvantages disadvantages disadvantages dis-
advantages disadvantages disadvantages disadvantages disadvantages disad-
vantages disadvantages disadvantages disadvantages disadvantages disadvan-
tages

info
info
info
info info info info info info info info info info info info

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages advantages
advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages

disadvantages disadvantages disadvantages disadvantages disadvantages dis-
advantages disadvantages disadvantages disadvantages disadvantages disad-
vantages disadvantages disadvantages disadvantages disadvantages disadvan-
tages

info info info info info info info info info info info info info info info info
info info info info info info info info info info info info info info info info
info info info info info info info info info info info info info info info info
info info info info info info info info info info info

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages advantages
advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages

disadvantages disadvantages disadvantages disadvantages disadvantages dis-
advantages disadvantages disadvantages disadvantages disadvantages disad-
vantages disadvantages disadvantages disadvantages disadvantages disadvan-
tages

info info info info info info info info info info info info info info info info
info info info info info info info info info info info info info info info info
info info info info info info info info info info info info info info info info
info info info info info info info info info info info info

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages advantages
advantages advantages advantages advantages advantages advantages advan-
tages advantages advantages advantages advantages advantages

disadvantages disadvantages disadvantages disadvantages disadvantages dis-
advantages disadvantages disadvantages disadvantages disadvantages disad-
vantages disadvantages disadvantages disadvantages disadvantages disadvan-
tages

info info info info info info info info info info info info info info info info info
info info info info info info info info info info info info info info info info info
info info info info info info info info info info info info info info info info info
info info info info info info info info info info info

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.12 Polyhedra

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.13 TimesTen

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.14 Csql

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.15 SolidDB

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.16 MonetDB

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.17 RDM Embedded

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.18 FastDB

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.19 QuiLogic

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.20 Pico4

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.21 Pico4v2

intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB
intorDB intorDB intorDB intorDB intorDB intorDB intorDB intorDB in-
torDB intorDB intorDB

2.22 Conclusion

Tabular result here!

Part II

Performance Analysis

Chapter 3

Performance Analysis

Introduction

“Every database vendor cites benchmark data that proves its database is the fastest, and each vendor chooses the benchmark test that presents its product in the most favorable light” [Burleson 02].

In this chapter we will introduce database performance analysis’ problem, talking about how to measure the performance, how to understand who is the fastest and how to choose the best database. This preface is not specific to in-memory databases, but can be generalized for every DBMS.

3.1 Impartiality Problem

The most important problem related to performance analysis, and benchmarking, is the validity, understood as impartiality, of the results.

3.1.1 Benchmark History

The first benchmarks were developed internally by each vendor to compare their database's performance against the competitors. But when these results were published, they weren't considered reliable, because there was an evident conflict of interest between the vendor and its database [Gray 93].

Benchmark Wars

Also when benchmarks were published by third parties, or even by other competitors, the results were always discredited by loser vendors, who complained for the numbers, starting often a benchmark war.

A benchmark war started when a loser of an important and visible benchmark reran it using specific enhancements for that particular test, making them get winning numbers. Then the opponent vendor again reran the test using better enhancements made by a "one-star" guru. And so on to "five-star" gurus.

Benchmarking

Benchmarking is a variation of benchmark wars. Due to domain-specific benchmarks there was always a benchmark which rated a particular system the best. Therefore the vendors promoted only benchmarks which highlighted the strengths of their product, trying to make them as a standard. This led to a proliferation of confuse benchmarks.

Although these phenomena were drastically reduced by the foundation of the Transaction Processing Performance Council (TPC) in 1988, they are

still alive.

3.1.2 Benchmark Differences

It is now evident how hard is to understand which database is the fastest. Benchmarks cannot be properly used to analyze the performance of several databases and choose simply the best. Every benchmark has a bias, testing some particular aspect of database systems, such as writes, reads, transactions and so on. Beyond, every benchmark operation can be implemented in different ways by the same database: creating a connection for every operation or using a connection pool; use different transaction isolation level; load the whole database in RAM or split it in different hard disk partition; etc. Furthermore the same benchmark, with the same implementation for every database, can show dissimilar results when it runs on different platforms (hardware and software).

All these differences are grouped by three categories:

1. Test scenario: reads, writes, etc.
2. Test implementation: there are different way to implement the same transaction.
3. Execution platform.

3.1.3 The Axiom

By this reasoning comes out the axiom which will guide the following chapters, and the work behind this thesis:

"There is not a slower or a faster database: databases are only slower or faster given a specific set of criteria in a given benchmark".

This sentence comes from an article of Bernt Johnsen [Bernt 05], who, in response to a benchmark comparing HSQL and Derby, stated that it's easy to make a benchmark, but it is always hard to communicate the meaning of the results. The interpretation depends on too many criteria, so that it's not easy to say which database is the fastest, unless specifying the set of criteria used in the benchmark.

3.2 Measures

While analyzing databases performance it's not possible to measure which database is the best, but we can measure many other parameters, which then can be interpreted in order to choose the database that fit our needs.

When choosing a database, the most important features to evaluate and compare are:

Throughput is the number of transactions per second a database can sustain. This is the most important feature to consider when evaluating a database system. The most representative application scenario to understand the meaning of the throughput is an on-line transaction processing system. This kind of application requires the database to sustain a certain number of transactions per second, based on the number of users, and not every database can suite the needs of the application itself. Another example is real time applications: they require even higher throughput. Therefore it is crucial to understand if a cer-

tain database is suitable to an application in terms of transactions per second.

Latency/responsiveness is the time the database takes to execute a transaction. It can be measured as the inverse of throughput.

File size of the database image. While traditional DBMS' store objects (tables), indexes and a transaction log file on the file system, in-memory databases, achieving durability, usually store only a journal file containing all the transaction executed on the database. Although this may seem a small file, it can become very huge, even more than the database image. This measure is interesting whereas each database use different data structures.

RAM usage is the quantity of RAM a process uses while running. Talking about in-memory database, this can be another way to measure the size of the database image. In addition, this is a critical value to take in mind: IMDBs only works correctly and efficiently under the hypothesis that the RAM is enough to contain the whole database.

CPU load becomes more important with the increasing number of services that live together in the same server. CPU is a precious resource shared by all the processes in a particular computer. While database systems can be usually deployed on a dedicated server, embedded databases live together at least with the application which use the database itself. And often in-memory databases are executed in embedded mode.

Disk I/O shows the usage of the hard disk, the bottleneck for every tra-

ditional databases. Although pure in-memory databases never access to the disk, when adding durability through a transaction log file, disk I/O is a bottleneck for IMDBs too. Therefore this measure is less important than the others because it is possible to take the hypothesis that every databases have the same disk I/O.

Startup time is the time the database needs to become operational.

Shutdown time is the time the database takes to shut down and kill the process.

3.3 Choosing a Database

From the previous sections, it's clear how it is difficult to use benchmarks to prove which database is the fastest. Even with a fair benchmark, which can be very useful to understand the performance of databases, it is still difficult to choose a database: performance is only one factor to consider when evaluating a database [Burleson 02].

Other factors to consider are:

- The availability of trained DBAs.
- The vendor's technical support.
- The cost of ownership.
- The hardware on which the database will be deployed.
- The operating system which will support the database.

In other words, it is very difficult to choose the right database for our needs, and, of course, while evaluating databases and benchmarks there is absolutely *no winner*.

Chapter 4

Test Suite

In this chapter we will define the tests used to run a benchmark and to analyze performance. The tests are divided into three categories: base test case, load test case and acid test case.

4.1 Why Define Test Scenarios

The axiom, previously described in paragraph 3.1.3 at page 36, expresses clearly the difficulty to analyze databases' performance and how every result obtained in a given benchmark depends on the specific set of criteria used in the benchmark itself.

In paragraph 3.1.2 there is also a description of the three major categories in which the criteria fall. The first of these categories is test scenario: different test scenarios may show completely different performance results. It's not possible to avoid this behavior, but we can define clearly every tests so that we will be aware of the differences between them.

4.2 Base Test Case

Base test case is a collection of very simple tests, which are configured mostly as a race. This is exactly what the major part of benchmarks does, particularly Poleposition [Grehan 05].

Every test can execute different read/write operations on the database and all these operations are inside a loop. The word *transaction*, used extensively in the following paragraphs, refers to an execution of the loop, and therefore all the operations inside it.

The key features of these kind of test are:

- A fixed number of transaction before the test stop: so tests are configured as a race where every database must execute a certain number of transaction.
- A fixed amount of time before the test stop: as an alternative to a fixed number of transaction, every test run for a specific amount of time, executing the maximum number of transactions per second.
- Different kind of object: tests must be able to create/retrieve/update/delete simple flat objects with few fields, or complex flat objects with many fields, or still hierarchical objects, and so on. This feature let us test effectively the performance on the objects used in the domain of our interest.
- Single task: to keep these test simple it's better avoid concurrent test, but it is still possible to implement concurrently the different operation executed on the database.

4.2.1 Race Test: Timed

We already said base tests are configured inherently as a race. These tests can be used to show the maximum throughput the database can reach doing a particular operation on a specific object. Metaphorically it's like a rocket car in the desert trying to reach its speed limit. In a real application this is rarely useful, but can give us an idea of database's limits.

In order to create a test scenario we have to define the object/s involved by the test, the operations on which the test loop and when the test should stop. For example:

- The object represent a person, with only two fields: an id number and a name. This is a simple flat object.
- The only operation executed by the test is a write operation: every time an object will be added to the database.
- 60 seconds is the stop condition of the test: after 60 seconds of execution of writes (objects person inserted in the database) the test stops.

Clearly the time is not a significant measure, every test's execution run for the same amount of time: 60 seconds. Instead of the time, a more interesting measure to take is the throughput. This test shows the maximum theoretical value for the throughput, which in a real usage scenario will never be outperformed.

This test, and its results, are not useful to understand the real performance and potentiality of the database and so to choose the database for our needs, but it can be used in an early stage to reduce the databases' set

which will be analyzed extensively with further tests. In other words we can throw away every databases whose maximum throughput is not enough to our needs.

4.2.2 Race Test: Transactions

This test is really similar to the previous test. The only difference is in the stop condition:

- The object represent a person, with only two fields: an id number and a name. This is a simple flat object (same as before).
- The only operation executed by the test is a write operation: every time an object will be added to the database (same as before).
- 1.000.000 of transactions is the stop condition of the test: when 1.000.000 of writes are executed (objects person inserted in the database) the test stops.

While in the previous test the duration was meaningless, this time, like in a race, it shows which database is the fastest (the winner of the race). Despite this, the duration is still of little use. Also for the throughput the same considerations made before are still valid.

But this new test is useful also to take other interesting measures, such as the file size, whose meaning has already been explained in paragraph 3.2. We can evaluate how the file size increase with the number of objects (person) inserted in the database. This is because, differently to the previous one, every test's execution perform a fixed amount of transactions.

4.3 Load Test Case

Restrictions introduced by base tests are very evident although they can be usefull in certain situation and for simple and fast approximate results. These restrictions are due to the impossibility to simulate real complex application scenario. Substantially the main restrictions are two:

1. Real application scenario are rarely single thread/task: to stick to the axiom at paragraph 3.1.3, the best way to understand which database fits our needs is to make test scnerios most realistic as possible.
2. The second restriction result from the first:trying to make test scenarios more realistic, it is necessary to introduce some sort of control for the throughput. When a test stress the database engine to its maximum level, some other mechanisms may not work properly, such as the garbage collector, caching, indexing, etc. In addition, when a test concurrently accesses the database, this restriction is even more evident: if a thread is not limited in its transactions per second, it may lower the performance of other threads.

Base test case offer very important and useful results, but when the need to test the average load of a real application of our interest, they are not enough. These restrictions are solved by load test case, which allow a better simulation of the application scenario. The key features of these tests, differently by base tests, are:

- Multi task: there will be not only a combination of read/write/update/delete in one transaction executed synchronizedly, but in more

transactions running concurrently against the database.

- A bond on the throughput: for every transaction it is possible to specify the maximum throughput, if the database can reach it, in order to simulate real application usage.

From these features comes the name "load test case", which means the simulation of an average, or specific, load on a certain database. This idea, and the need of this kind of tests, can simply be explained by a Bernt's metaphor, who has been already cited for the axiom at paragraph 3.1.3. The metaphor is:

"I don't buy the fastest car in the market. I don't even buy the fastest car I can afford. I buy a car I can afford that fits my needs... e.g. drive to the northern Norway with 2 grown-ups, 3 kids, a kayak, sleeping bags, a large tent, glacier climbing equipment, food, clothes and so on. Can't do that with a Lamborghini" [Bernt 05].

This metaphor explain exactly the need not for the fastest database, but for a database which can sustain the load of the application without any complexity during its normal functioning, such as a database snapshot freezing all writes operation, or a bug/memory leak making the database crash.

4.3.1 Real Time Prepaid System

Keeping these features in mind, we want to create a load test case to analyze exactly which performance a database system offers for a practical use case: a real time prepaid system, such as a telephone company. Basically

the test is the concurrent execution of 3 different task, each involving complex transactions: check balance, accounts handling, and manage calls and services.

We start describing the domain objects used by this test, and then we pass to analyze the three different tasks involved.

Domain Objects

Domain objects involved by this test are typical for a telephone company who handle telephone calls for every person, which is represented by an account, and which can access to many services through a web application, offered by the company for its customers. There are four domain objects used by this test:

Account : this object represents a customer in the real time prepaid system.

It contains all the typical information needed by a telephone company, such as the balance, the type of subscription, etc. In other words, this is a complex flat object, that is having many private fields but no hierarchies. In a real application there are millions of accounts instantiated in the database, one for every customer. This dimension is also very important to make the simulation as real as possible, in order to get realistic results.

Msisdn : it is the unique number which identifies a subscription in a GSM or UMTS mobile network, it is the telephone number. Each msisdn object is linked to an account. Therefore there are also millions of msisdn objects in the database, referring to a real application.

Webuser : this represents the customer logged in the company web site.

It contains the username and the password to access to web services.

In common with msisdn, it is linked to an account too, and both are simple flat objects: a very simple object with few fields.

Session : when a customer start a call, an object session is created, and it keep all the information about the call which is going on, until the call ends. After the call this object doesn't is deleted. A session contain the start time of the call represented, and the time of the last unit. Each session references to an account, which started the call. But there are not as many sessions as accounts, not everyone will start a call in the same time, except New Year's Day!

Check Balance Task

This task simulate a customer checking his residual balance for the prepaid card, a SIM in the case of the telephone company. First, the customer logins in the website or calls the dedicated number, and then read/listen his balance.

Each task, as already explained, corresponds to a transaction composed by different operations. This leads to illustrate how the task work in terms of operations executed:

1. *Read* randomly the msisdn if the customer makes a call or the webuser if he access to the website.
2. *Read* the account referenced by the msisdn or webuser, and then check the residual balance, a simple account's field.

So this task executes a total of two read for every transaction.

Another important parameter to make this task a part of a load test is the amount of transactions per second this task should sustain. The average throughput for check balance task, considering there are millions of accounts in the database, is about ten transactions per second.

Accounts Handling Task

The account handling task simulates the creation of a new subscription by a customer, and so the creation of the account object, and after the relative msisdn and webuser.

The operations involved by this task are:

1. *Write* of a new account: all informations of a customer are inserted in the system by creating a new account object.
2. *Write* of a new msisdn, containing the telephone number of the new subscription.
3. *Write* of a new webuser with username and password for the customer.

To sum up, this task executes three write on the database for every transaction.

To simulate a real scenario, the average throughput is about ten transactions per second.

Manage Calls and Services Task

This task simulates a call started by a customer. After checking the balance, and therefore if the account is allowed to start a call or use a service, a

new session is created and updated every 30 seconds, the unit time. On the average a call lasts about two minutes. After this time the session is deleted.

In terms of operations executed on the database, it can be described as follow:

1. *Read* the msisdn of the customer starting the call.
2. *Read* the account referenced by the msisdn, and its residual balance, and the other parameters, to check the user's permissions for the service requested.
3. *Write* a new object session and *update* account's balance.
4. *Update* the session every 30 seconds and the account's balance.
5. *Delete* the session after the call is ended.

This is the most complex and important task, because it represents a very frequent task. No wonder if the average throughput is about 2000 transactions per second.

Real Time Prepaid System Summary

To sum up this whole load test we have to define the objects used by the test, the task concurrently executed, and when the test stop.

There are four domain objects involved by this test. Three of these are simple flat objects: msisdn, webuser and session. The forth is a complex flat object: account. The number of objects already initialized in the database, before the test start, is on the order of millions.

Three task are executed concurrently. One is the major task, simulating a customer's call, and it runs 2000 times per second. This task is the bottleneck of a real application, and it could be also tested alone in a base test to understand the database's limits in running this task, and therefore to have an idea of the database potentiality. The other two are minor task, in fact the throughput is limited to ten transactions per second.

This test is a load test and therefore we are not interested in stopping it after a certain amount of transactions executed. But what we want is to let the test run for many minutes until to many hours to understand if the database can sustain the load generated by the test. Then the stop condition will be a specific amount of time.

4.4 Acid Test Case

Base and load test case are used to analyze databases' performance, but, as already said, a database is not made only by performance: there are many other parameters to take into consideration before judging a database, such as ACID properties. To make these tests an all-round test suite, we could add some tests for verifying databases' ACID properties. The above corresponds to the work done by the Transaction Processing Performance Council with their TPC-C benchmark [Council 07]: performance is not the only benchmark's goal, but ACID properties are tested too.

Especially the durability property is most interesting to analyze when talking about in-memory databases. Pure IMDB have no durability at all, while it can be added in different degree, as already described in paragraph

1.5. Therefore would be very useful to know, except the vendor's words, how strong is the durability solution implemented by the in-memory database.

Nevertheless, testing ACID properties is not a simple goal to achieve. Particularly the durability is the hardest property to be tested. Also the tests developed by TPC are not intended to be an exhaustive quality assurance test.

Chapter 5

Database Benchmark Softwares Overview

5.1 Benchmark Requirements

Desirable attributes ... list of attributes [DeWitt 97].

5.2 The Open Source Database Benchmark

5.3 Transaction Processing Performance Council

*Benchmark results are highly dependent upon workload, specific application requirements, and systems design and implementation. Relative system performance will vary as a result of these and other factors. Therefore, TPC-C should not be used as a substitute for a **specific customer applica-***

tion benchmarking when critical capacity planning and/or product evaluation decisions are contemplated [Council 07].

5.4 Apache JMeter

5.5 Poleposition

Poleposition fa schifo! [Grehan 05]!!!

Chapter 6

The New In-Memory Database Benchmark Application

6.1 Functional View

6.2 Development View

6.3 Plug-In Architecture

Chapter 7

Results' Analysis

Appendix A

The first appendix

questa e un appendice

Bibliography

- [Bernt 05] Johnsen Bernt. *Database Benchmarks*. http://weblogs.java.net/blog/bernt/archive/2005/10/database_benchm_1.html, october 2005.
- [Burleson 02] Donald Burleson. *Database benchmark wars: What you need to know*. http://articles.techrepublic.com.com/5100-10878_11-5031718.html, august 2002.
- [Council 07] Transaction Processing Performance Council. *TPC Benchmark C*. www.tpc.org/tpcc/spec/tpcc_current.pdf, june 2007.
- [DeWitt 97] Dave DeWitt. *Standard Benchmarks for Database Systems*. <http://www.tpc.org/information/sessions/sigmod/sigmod97.ppt>, 1997.
- [Fowler 05] Martin Fowler. *InMemoryTestDatabase*. <http://martinfowler.com/bliki/InMemoryTestDatabase.html>, november 2005.

- [Gorine 04] Andrei Gorine. *Building durability into data management for real-time systems*. <http://www.mcobject.com/downloads/bassep04p19.pdf>, september 2004.
- [Graves 02] Steve Graves. *In-Memory Database Systems*. <http://www.linuxjournal.com/article/6133>, 2002.
- [Gray 93] Jim Gray. *Database and Transaction Processing Performance Handbook*. <http://research.microsoft.com/~Gray/BenchmarkHandbook/chapter1.pdf>, 1993.
- [Grehan 05] Rick Grehan. *An Open Source Database Benchmark*. <http://today.java.net/pub/a/today/2005/06/14/poleposition.html>, june 2005.
- [Hobbs 03] Darren Hobbs. *Prevayler Pertubations*. <http://darrenhobbs.com/?p=225>, march 2003.
- [Miller 03] Charles Miller. *Prevayling Stupidity*. http://fishbowl.pastiche.org/2003/04/11/prevayling_stupidity/, april 2003.
- [Prevayler 08] Prevayler. *Prevayler Home Page*. <http://www.prevayler.org/wiki/>, march 2008.
- [Wuestefeld 01] Klaus Wuestefeld. *Object Prevalence*. <http://www.advogato.org/article/398.html>, december 2001.