

# XML and NoSQL DBMS: Migration and Benchmarking

Author:  
Prakash Thapa

University of Konstanz

June 22, 2015

## **Abstract**

As a dynamic data is growing rapidly, database vendors search for the alternative of the classical database management system such as RDBMS. XML and NoSQL databases are two non-relational database management systems growing since the popularity of Big Data. This thesis focus on the comparative analysis of these two new database systems based on use cases and existing solutions. On the first part, our focus will be data migration from XML database to NoSQL database. We will also discuss the performance of both systems based on the standard queries.

.....

## **Zusammenfassung(German Abstract)**

XML und NoSQL

## **Acknowledgments**

.....

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	Contribution . . . . .	2
1.3	Overview . . . . .	2
<b>2</b>	<b>NoSQL and XML Databases</b>	<b>3</b>
2.1	XML Database . . . . .	3
2.1.1	BaseX . . . . .	3
2.2	NoSQL Database . . . . .	3
2.2.1	MongoDB . . . . .	5
	Data design . . . . .	5
	Embedded . . . . .	5
	Reference . . . . .	5
	Document . . . . .	6
	Data-types for Document Fields . . . . .	6
	Indexing . . . . .	7
	Query Model . . . . .	7
	System Architecture . . . . .	7
2.2.2	Couchbase Server . . . . .	8
	Metadata . . . . .	8
	Document key . . . . .	9
	Bucket and vBucket . . . . .	9
	Data Model . . . . .	9
	Querying and Indexing . . . . .	9
2.2.3	RethinkDB . . . . .	10
	Data Model . . . . .	11
	Query Model . . . . .	11
	Indexing . . . . .	11
<b>3</b>	<b>Semi-structured Data: XML and JSON</b>	<b>13</b>
3.1	Problem to translate from XML to JSON . . . . .	13
	Root node and anonymous values . . . . .	13
	Arrays . . . . .	14
	Identifiers . . . . .	14
	Attributes . . . . .	14
	Namespaces . . . . .	14
	Others . . . . .	14
3.2	Mapping . . . . .	14
3.3	Migration from XML to JSON . . . . .	15
	Step 1 . . . . .	15

Step 3 . . . . .	16
3.4 XMark . . . . .	18
3.4.1 Dataset . . . . .	18
3.4.2 Queries . . . . .	19
3.4.3 XMark data into NoSQL Database . . . . .	21
3.4.3.1 XMARK in MongoDB . . . . .	22
3.4.3.2 Couchbase Server . . . . .	22
Document design . . . . .	22
3.4.3.3 Rethinkdb . . . . .	23
<b>4 Benchmarking</b>	<b>25</b>
4.1 Benchmarking configurations . . . . .	25
4.2 Results and Analysis . . . . .	25
4.2.1 Summary . . . . .	26
4.3 Discussion . . . . .	26
4.4 Conclusion . . . . .	26
4.5 Future Work . . . . .	26
<b>Appendices</b>	<b>29</b>
<b>A XMARK Queries</b>	<b>31</b>
A.1 MongoDB . . . . .	31
A.2 RethinkDB . . . . .	35
A.3 Couchbase . . . . .	39

# Chapter 1

## Introduction

### 1.1 Motivation

As the digital world growing very fast, the massive amount of data collected today in the field varying from business to scientific research, is becoming more complex for storage, querying and providing mechanism for failover. As the data collection grows so largely, the traditional data management tools e.g. relational database management systems(RDBMS) are struggling to handle such data effectively. The pre-design rigid schema structure of RDBMS made more complex for variety of data. High data velocity where massive read and write operations are possible from different geographic locations, storage of structured, semi-structured and unstructured data in gigantic volume, together is referred as *Big Data*. This has become a global phenomena which added more difficulties on RDBMS. RDBMS are not able to store and process big data (that commonly has unstructured data) effectively and are not very efficient to make transactions and join operations (Abramova and Bernardino, 2013).

New database technologies NoSQL and XML Databases came into existence to overcome above mentioned problems of RDBMS. These DBMS store and process unstructured and semi-structured data(data that do not follow specific data format varies and can grow in real time). These systems are not just able to solve the issues of existing data management systems, but also rise as alternative in their way. Both these new technologies focus on varieties of data in large volume, dynamic schema and scalability.

There are some research analysis between NoSQL and RDBMS. Nance et al. (2013) examine the pros and cons of both NoSQL and RDBMS, Cattell (2011) analyze similarities between scalable SQL and NoSQL databases, Hadjigeorgiou et al. (2013) compare the performance between MySQL cluster as RDBMS and MongoDB as NoSQL Database. There are also some research related to RDBMS with XML or XML Database( Jiang et al. (2002), Shanmugasundaram et al. (1999)). But there are not much more research in the field of NoSQL and XML database together. So this thesis focus on these two new database systems. Migration of Data from XML database to NoSQL as well as performance of both system based on standard query will be examined.

## 1.2 Contribution

The main contribution of this thesis is that it provides the necessary techniques and algorithms for migrating data from an XML database to NoSQL databases. More specifically, it will focus on databases like MongoDB, Couchbase and RethinkDB as NoSQL databases and BaseX as XML Database. To approach this challenge, it is first necessary to understand the general architecture and data model of each of these databases as well as the way how they are queried. The performance comparison of these two systems will be based on XMARK dataset and its 20 standard queries (Schmidt et al., 2002). These 20 queries of XMARK project will be translated to each of NoSQL databases.

## 1.3 Overview

This master thesis is organized as follows: In chapter 2, we will introduce the two new databases systems. Chapter 3 defines the techniques and necessary algorithms to convert XML to JSON. In same section, the XMARK dataset will be introduced and procedure to store in NoSQL databases. In Chapter 4, the spotlight is on the performance tests and comparative analysis of each of the NoSQL databases with BaseX, an native XML database, based on the XMark benchmarking project.

# Chapter 2

## NoSQL and XML Databases

### 2.1 XML Database

Native XML Databases (NXDs) are document oriented databases that store XML data and all components of XML based on a specified logical model. This model might differ to individual NXDs, but all databases include at least elements, attributes, PCDATA, and document order. An XML document is the fundamental unit of storage, So each document should be in a valid XML format. XQuery 1.0 and XQuery 3.0 are standards query languages recommended by W3C for NXDs. XQuery includes XPath as sub-language which is a language for addressing specific part of an XML document. In addition to XQuery and XPath, XML databases also support XSLT, a method of transferring XML to other documents like HTML, plain text and XSL Formatting Objects. An extension to the XQuery/XPath specifies the full-text search queries ???. XML databases accommodate the Web standards of of the World Wide Web Consortium(W3C) including XML Schema, XPath, XQuery, Full-Text Search, XProc, XForm and XSL-FO. BaseX, MarkLogic, eXist-db, Sedna are some of representative example of NXDs.

#### 2.1.1 BaseX

BaseX is a Native XML database management system and XQuery Processor. It uses tabular representation of XML data tree to store XML document (Grün, 2010). In this dissertation, BaseX is represented as XML database. ....

### 2.2 NoSQL Database

NoSQL database "Not Only SQL" is distributed data management system for large and variable data. The main focus of any NoSQL database are scalability, performance and availability (Hecht and Jablonski, 2011).

*Scalability* is one of the basic goals of the NoSQL database system to store unstructured and semi-structured data over a distributed environment where data instances are large and stored separately in different nodes. These databases are designed to provide a facility for storing very large data capacity( from terabytes to petabytes)Abramova and Bernardino (2013). To achieve this objective, NoSQL databases mainly trade-off its capacity to handle ACID transactions to high scalability over a large number (hundred to thousands) of nodes. Since the data is stored in commodity hardware throughout many servers, network failures are common when

processing data operation. If one or more nodes unable to deliver the request of a client, the other nodes in the system can complete the operation without any loss. Data is replicated throughout different nodes for fault tolerance and systems allows eventual consistency among the replicas of same the piece of data Gao (2013). Based on the Data model and design architecture, NoSQL databases are categorized into 4 groups:

- **Key/Value store** The data representation of key-value stores are based on attribute pair, the data model is expressed as collection of  $\langle key, value \rangle$  tuples. The key is unique and value can be varied. The query operation on data can be uniquely performed by a key. There are no alternative ways to access or modify data without key. DynamoDB, Riak, Redis are some examples of key-value stores.
- **Document oriented databases** are based on semi-structured data model, where unique key stores a value that contains a tree like structure called *document*. Each document is a record with potentially varied attributes. These attributes are key-value pair in JSON or in JSON like data format where data can be accessed by using key or specific value (Hecht and Jablonski, 2011). MongoDB, Couchbase, RethinkDB, and CouchDB are some of representative systems in this category.
- **Column family databases** are also known as wide-column databases, which store data tables as a section of a column. Each column has its own index structure. These systems are suitable for heavy write and optimized for writing data to a smaller subset of records. Cassandra, BigTable, HBase are categories as wide column database. Figure 2.1 shows a column family record.
- **Graph databases** are based on graph theory. In these databases, data are represented as a graph in the form of nodes and edges similar to a social network. They store entities and relationships between these entities. Neo4J and OrientDB are two example of these databases.

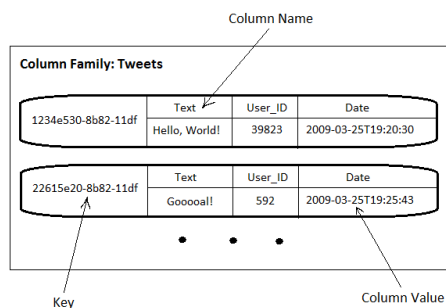


Figure 2.1: column based database

```
{
  "_id": "4c6cd19abac7061798000002",
  "CityId": 42231,
  "CountryID": 1,
  "RegionID": 833,
  "City": "Herat",
  "Latitude": 34.3330001831055,
  "Longitude": 62.2000007629395,
  "TimeZone": "+04:30",
  "DmaId": 0,
  "County": "HERA",
  "Code": ""
}
```

Figure 2.2: A document in document based database(JSON)

As there are many categories and varieties of database system, our focus in this thesis is limited to document oriented databases of NoSQL databases.



### 2.2.1 MongoDB

MongoDB is a schemaless document oriented database developed by MongoDB Inc.( then 10gen Inc.) and Open Source Community. It is intended to be a flexible, fast and Multi-datacenter scalable system.

**Data design** MongoDB stores data as documents with unique key. Beside querying a document using primary key, it allows to query by every attribute of document. This made possible to index every field of a document. Documents with similar schema are grouped together in a collection. A database contains one or more collections. As MongoDB is a schema-free system, a collection may contain any type of documents but mostly it has documents with similar schema. Collections do not enforce document to structure data rather as per requirements of the our application. Documents are modeled as BSON format which is a binary variant of JSON that supports additional data types like ObjectId, timestamp, datetime etc. In MongoDB, there are two principles that allow application to represent documents and their relationship: *reference* and *embedded* documents.

**Embedded** Embedded documents captures relationships between the data by storing related data in a single document structure. The documents in this method are structured as sub-documents in the form of Array or/and Object (Gao, 2013).

**Reference** MongoDB has no support for joins, So related data are stored in a single document. In some cases, it makes sense to store related data in separate documents, typically in different collection or even a database. Reference stores the relationships between data by including links and references from one document to another as described in Figure 2.3a. These references can be created in two way.

- **Manual reference** The application handle the relationship between the related data in manual reference. The `_id` field of one document is referenced to another document.
- **Database reference(*DBRefs*)** The *DBRefs* references one document to another documents using first document's primary index field, collection name and optionally database name. With the reference of collection, a document can be reference to another document belonging different collection where as referencing both collection name and database, MongoDB allows a document reference to multiple database.

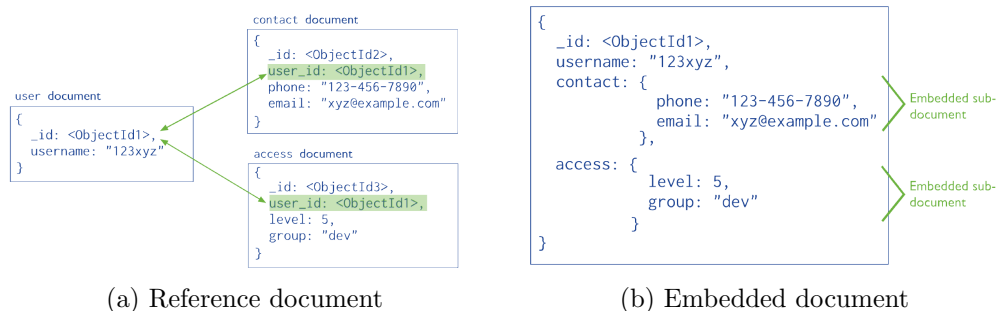


Figure 2.3: MongoDB document structure

**Document** A document is an abstract and storable unit in MongoDB. It is the data structure expressed in the form of JSON. All the accessible data including database record, query selectors, index specifications, server configurations etc. are represented as documents. The attributes name(field name) of a document can be a valid JSON string with three exception. *\_id* should be unique as it is primary index, a field name cannot start with dollar (\$) character, dot(.) character and finally, field name cannot have a null character.

```
{
    _id : "001"
    title : " MongoDB ",
    last_editor : "172.5.123.91" ,
    last_modified : new Date ("9/01/2015") ,
    body : " MongoDB is a..." ,
    categories : [" Database ", " NoSQL ", " Document
        Database " ] ,
    reviewed : false
}
```

Code 2.1: MongoDB sample document

**Data-types for Document Fields** MongoDB support following datatypes

- scalar types: **boolean**, **integer**, **double**
- character sequence type: **string**(encoded in UTF-8), **regular expression**
- **Object**: All BSON-Objects
- **ObjectId** : It is BSON type 12 byte long binary value to uniquely identify the document. Documents are stored in a collection with *\_id* field which also act as primary key with the value of ObjectId. The object id datatype is composed of the following components:
  - timestamp in seconds since epoch (first 4 bytes)
  - id of the machine assigning the object id value (next 3 bytes)
  - id of the MongoDB process (next 2 bytes)
  - counter (last 3 bytes)
- **null**
- **array**
- **date**

**Indexing** Each document in MongoDB is uniquely identified by a field *\_id* which is a primary index. Hence, the collection is sorted by *\_id* (Gao, 2013). In addition to primary index, MongoDB provides the mechanism to create secondary indexes for all the document fields. It supports various user defined indexes for field values of documents including single field index, multikey index, multidimensional index, geospatial index, text index and hash index. Single field, multidimensional and multikey indexes are organized using B-tree, whereas geospatial index is implemented using quad trees.

- *Single field index* only includes data from a single field of documents in a collection.
- *Compound index* holds reference to the multiple fields within a collection's documents.
- To index a field that contains an array value, MongoDB provides special indexing called *Multikey index*.
- *Text index* helps efficient search of a string in documents.

**Query Model** Queries in MongoDB are expressed in a JSON like syntax and send to MongoDB as BSON objects by database drivers(Orend, 2010). MongoDB's query can be implemented over all documents inside collections, including embedded object and arrays. The Query model support following features:

1. Queries over documents, embedded subdocuments and arrays
2. Comparison operators
3. Conditional Operators
4. Logical Operators: AND and OR
5. Sorting
6. Group by
7. Aggregation per query

In addition to this, MongoDB provides features for complex aggregation with the use of MapReduce. The result of MapReduce either can be stored as a collection or returns value to the client(Orend, 2010).

**System Architecture** MongoDB can be run in two modes. In stand-alone mode, a single *mongod* demon is running in a single node without any distribution. The other mode is sharded mode, where various services of MongoDB are distributed to several nodes.

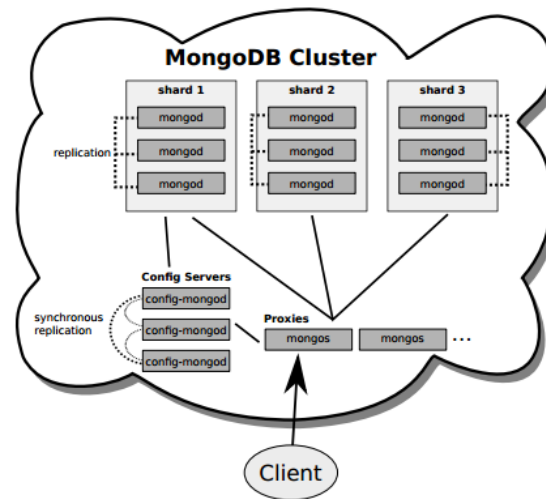


Figure 2.4: MongoDB Clusters

### 2.2.2 Couchbase Server

Couchbase Server is designed to be an operational data store which means it is solution for real-time data access. It is NoSQL database that facilitates both as a key-value store as well as document store. As key-value store, it is able to store multiple data type such as strings, numbers, datetime, and booleans as well as arbitrary binary data. The key-value generally treated as opaque Binary Large Object(BLOB) and do not try to parse it at the time of query. For document store, data need to be stored in the valid JSON format. Data in Couchbase Server are stored in logical unit called Buckets. Buckets are isolated to each other which have their own RAM quota and replica settings. These buckets can be technically compared as **database** in MongoDB or other RDBMS. Couchbase recommends as few as possible number of buckets even it is possible to have up to ten buckets in a single cluster. These buckets theoretically contain any type of data. All data type other than JSON can be retrieved only by their key. So, it is important to check meta type of data that are stored in a single document before retrieval.

**Metadata** For every value stored in a bucket, Couchbase Server generate following meta information that is associated with the document (Ostrovsky and Rodenski, 2014, p. 26).

- **Expiration** The Time to Live(TTL) also named expiration time is the life time of the document. By default it is 0 which indicates it will never expire, also can be set as Unix epoch time after which document is removed.
- **Check-and-Set(CAS)** The CAS value is 64-bit integer that is updated by server when associated item is modified. It enables to update information only if unique identifier matches the identifier of the document that needed to be update. CAS is used to manage concurrency when multiple client attempts to update the same document simultaneously.
- **Flags** are 32 bit integer and are set of SDK specific use for SDK specific need. For example, format in which data is serialize or data type of the object being stored.

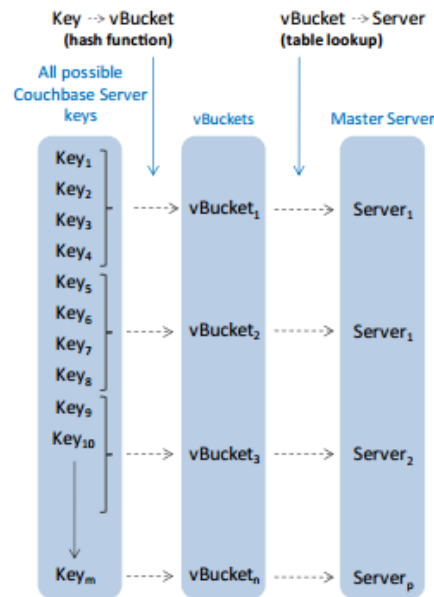


Figure 2.5: Couchbase Server vBucket

In addition of expiration, CAS and flags, three other meta information are stored at the time of document creation. The document's key *id* is saved as part of metadata, *type* is the type of a document either **json** for all valid JSON documents or Base64 encoded string for all other than JSON.

**Document key** Every value in Couchbase Server is saved in unique key called document key. Unlike MongoDB, Couchbase do not generate document key automatically, it should be created a string value upto 250 characters manually.

**Bucket and vBucket** Couchbase Server uses data bucket as a logical container of information. It provides a logical grouping of physical resources within a cluster (Lichtenberg and Ward, 2013). Contrast to MongoDB, the Couchbase Server do not have the concept of collections. Documents do not have their own fixed schema, multiple documents with different schema can be stored in same bucket. One or more attributes are added to differentiate the various objects stored in a bucket and create indexes on them.

Each bucket is split into 1024 logical partition called vBuckets. A vBucket is treated as owner of subset of key. Every key belongs to a vBucket.

**Data Model** A document is the most basic unit of data manipulation in Couchbase Server. All the documents are stored in JSON format without predefined schema.

**Querying and Indexing** Couchbase Server's views are responsible for querying data from bucket. Views are defined in a specific kind of document called *design document*. These documents are bounded to a single bucket, hence cannot executes from others. The design document holds JavaScript code that implements *Mapreduce* operations which creates view's index in user defined format. The *map()* function in design document filters and extracts information from multiple documents that are stored in a bucket. For a view, each document in a bucket is submitted to the *map()*

function at once. After filtration, the output is a set of key/value pairs which are the results of a view. The *map()* function are supplied with two parameters. First, the JSON document that is stored in a bucket and optional second parameter is metadata associated with the document. The *map()* function's output can be zero or more "rows" according to the filter used in *emit()* function. Each *emit()* function returns a single row but can be called multiple times inside a single map function. The *reduce()* function is used to aggregate the content generated in map the phase. Couchbase has built-in reduce functions like *\_count*, *\_sum* and *\_stats*. A typical example of Couchbase Server's MapReduce is illustrated in Table 2.1

Table 2.1: Mapreduce in Couchbase Server

<i>map()</i>	<i>reduce()</i>
<pre> function (doc, meta) {   if(doc.doctype &amp;&amp; doc.doctype=="     closed_auctions"){     if(doc.price){       if(doc.price &gt;= 40) {         emit(doc.id,doc.price)       }     }   } } </pre>	<b>_count</b>

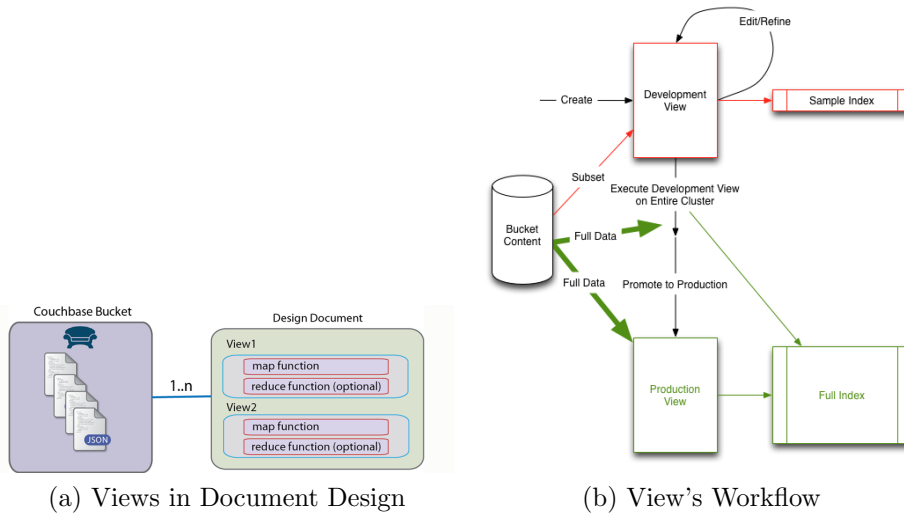


Figure 2.6: Couchbase Server's Document Design

### 2.2.3 RethinkDB

RethinkDB (Rethinkdb, 2015 (accessed January 12, 2015)) is a distributed database system to store JSON documents. It has query languages named ReQL which automatically parallelize queries in multiple machines. ReQL is based on three main principle. Firstly, it is completely embedded with programming language.

Secondly, ReQL queries can be passed as pipeline from one stage to another to get required result that means it is possible to use series of simple queries together to perform complex operation. Finally, all the queries are executed in server without any intermediate network round trip required between the server and clients.

**Data Model** Rethinkdb stores JSON documents with binary on disk serialization. The data types supported by RethinkDB is same as of JSON. There are two ways to model relationship between the documents in RethinkDB.

1. In *Embedded arrays*, the related sub-documents are inserted with a specific key inside a document as in MongoDB 2.2.1. The advantages of using embedded arrays are:
  - The Queries are tend to be simpler.
  - If a dataset is large and don't fit in RAM, data is loaded faster from disk as compare to tables.

Disadvantages of embedded arrays are:

- Everything should be loaded in memory for any operation.
2. In *multiple tables* approach, document with similar schema are stored in tables which are connected by reference key to another table, like in relational database. Unlike embedded document, operation of a table don't required to load data from reference table. Table has also some disadvantages like linking between table is more complicated.

**Query Model** ReQL is the RethinkDB's query language. In ReQL, the JavaScript expression can be embedded at anywhere as a part of query as it embeds into programming language. The lamda function is a part of query or sub-queries that gives more flexibility for expected result. All the ReQL queries are chainable, at any part of query the with dot `{.}` operator can be extended at any level as in Code 2.4.

```
r.table("users").run(conn)
r.table("users").pluck("last_name").run(conn)
r.table("users").pluck("last_name").distinct().run(conn)
r.table("users").pluck("last_name").distinct().count().run(conn)
```

Code 2.4: Chainable Query in ReQL

Unlike other NoSQL databases, it supports join queries like in relational database in one to many or many to many relationship based on requirement of application in distributed manner.

**Indexing** RethinkDB uses B-Tree to store indexes. At the time of table creation, it provides a option to specify the attribute as a primary key. By default, *id* behave as a primary key if nothing is specified. At the of time data insertion, if specified primary key attribute exists, it's value is used to index the document. Otherwise, a random unique string is generated to index automatically. Beside primary index, RethinkDB support compound, secondary, geospatial and arbitrarily computed indexes. Every RethinkDB query, including update operations, uses one and only one index(Rethinkdb, 2015 (accessed January 12, 2015)).

- Simple indexes are created to get or sort efficiently by a value of an attribute.

- Compound indexes as name suggest, are used to retrieve document by multiple fields.
- Multi indexes are the indexes for array value of a field.
- Any arbitrary expression can be used to create index at any type of user defined expressions like lamda functions.



# Chapter 3

## Semi-structured Data: XML and JSON

XML, by definition a textual markup language where data elements are ordered by nature: *string* is core data type from which richer data types e.g. integers, floats, user-defined abstract data and so on can be derived (Schmidt et al., 2002).

A JSON or JavaScript Object Notation is a programming language model. It consist of minimal textual representation and a subset of JavaScript. JSON document consists of two data structures:

- Object, an unordered collection of key-value pair that are encapsulated by curly braces { and }. The key is a string encapsulated in double quotes and must followed by a column(:) to it's value. The key must be unique for each object.
- Arrays, which are an ordered list of values

A JSON value can be Object, Array, number, string, true, false and null.

### 3.1 Problem to translate from XML to JSON

JSON and XML look conceptually similar as they both are text based markup language, which are designed to represent data in human-readable form, exchangeable across multiple platforms and parseable by common programming language. When we look at them first, they appear to be quite similar, with difference only in their syntax (Lee, 2011). But it turns out that they are fundamentally incompatible, as listed below:

itemize followi  
tent

**Root node and anonymous values** Each XML document have one and only one root node. JSON support the anonymous values(also refer as string value) that don't need key-value pair. For example, Table 3.1 the XML root node is implicitly created in the model and will not have a textual representation in JSON.

Table 3.1: Anonymous values of JSON in XML

JSON		XML
		<root>Hello World</root>
["Hello World"]	or	<root value="Hello World">

**Arrays** Arrays are native data types of JSON which does not exists in XML. There is no direct markup for array of Object from JSON to XML.

**Identifiers** XML is much more restrictive for identifiers as compare to JSON which allows any string to be an identifier. Translating from XML to JSON do not cause problem but in reverse case might lead to invalid XML. For example "Hello World" is a valid identifier in JSON but not valid attribute or element in XML as there exists a whitespace between these two words.

**Attributes** JSON does not have the concept or representation for XML's attributes. When mapping from XML to JSON, attributes are translated to name object member along with other child elements. This information will be loss when mapping from JSON to XML.

**Namespaces** XML supports the concept of namespaces to provide uniqueness in element and attributes in XML document. Namespaces do not exists in JSON. Mapping QNames in XML with namespaces to JSON can lead to ambiguous and duplicate names.

**Others** There are also some problem like processing instructions and comments which XML supports but not JSON. Other issues for example character set and encoding are not easily exchangeable.

## 3.2 Mapping

XML and JSON have different data types. Compare to JSON, XML has more flexible data types. In Table 3.2 is given types of data for simple XML and relevant JSON type.

Table 3.2: Translation of simple XML data types into JSON

XML type definition	JSON type definition
xs:string	{ "type": "string" }
xs:boolean	{ "type": "boolean" }
xs:float xs:double xs:decimal xs:integer (All Other Numbers)	{ "type": "number" }
(Remaining all others)	{ "type": "string" }

For complex XML data type, there are two options in JSON, either object or array based on XML object.

### 3.3 Migration from XML to JSON

It takes different stages to convert XML data to JSON.

Step1.

XML to JSON friendly XML: Convert standard XML data to JSON friendly XML.

Step2.

Data type of XML to JSON data type: Define the data type of an XML element whether it is array type or object type or other scaler type.

Step3.

Final steps: Map XML to JSON object or array.

#### Step 1

Algorithm 1 provides psudo-code for Step 1. At first, all attributes of XML document are represented as ordered(first) child element then attributes are deleted. There can be more than one attributes which are inserted in order. XML node can contain both text node and element node as a child, but JSON can have only key-value pair or value only for array. So each of text node which also has element node as siblings moved to "childtext"(name of the element) node. At the end of first step, if an element has another element as a child node, then it is represented as object type of JSON. It is necessary to identify if a siblings of an element node( $S$ ) has same element name or not. If this condition exists,  $S$  is represented as the array in JSON. All the child nodes of all siblings are moved inside  $S$  node. If  $S$  is already object type then it is replace with array type else new array is added.

---

#### Algorithm 1: Pseudocode to convert normal XML to JSON friendly XML

---

```

Initialize  $D = \text{"XMLdocument"}$ ;
for all descendant-or-self node of  $D$ ,  $X$  which has attributes  $A$  do
    move all attributes  $A$  to ordered child element of  $X$ 
end for
for all descendant-or-self node of  $D, X$  do
    if  $X$  contains Text node and element node Both then
        create new element "childtext"
        move text node to "childtext" element
    end if
end for
for all descendant node of  $D, X$  do
    for all child element  $C$  in  $X$  do
        if  $C$  has siblings  $S$  with same name then
            convert  $C$  as Array type
            move child of  $S$  into  $C$  as  $\langle \_ \rangle$  element
        end if
        if  $C$  has child element then
            convert  $C$  as Object type
        end if
    end for
end for

```

---

**Algorithm 2:** convert data type of XML to JSON data type(Step 2)

---

```

Initialize  $D = \text{"XMLdocument"}$ ;
for all descendant-or-self node of  $D$  as  $C$  do
  if  $C$  has no attribute "type" then
    get content of  $C$ 
    identify content type according to Table 3.2. and add attribute to  $C$ 
    "type"
  end if
end for

```

---

Table 3.3: XML to JSON friendly XML(STEP 1)

XML	Algorithm 1
<pre> &lt;info&gt;   &lt;name&gt;     &lt;f&gt;a&lt;/f&gt;   &lt;/name&gt;   &lt;age&gt;24&lt;/age&gt;   &lt;ismarried&gt;&gt;false&lt;/ismarried&gt;   &lt;city name="Armonk"/&gt;   &lt;state&gt;NY&lt;/state&gt;   &lt;contact&gt;     home     &lt;p&gt;993-330&lt;/p&gt;     &lt;p&gt;993-331&lt;/p&gt;   &lt;/contact&gt; &lt;/info&gt; </pre>	<pre> &lt;info type="object"&gt;   &lt;name type="object"&gt;     &lt;f&gt;a&lt;/f&gt;   &lt;/name&gt;   &lt;age&gt;24&lt;/age&gt;   &lt;ismarried&gt;&gt;false&lt;/ismarried&gt;   &lt;city type="object"&gt;     &lt;name&gt;Armonk&lt;/name&gt;   &lt;/city&gt;   &lt;state&gt;NY&lt;/state&gt;   &lt;contact type="object"&gt;     &lt;childtext&gt;home&lt;/childtext&gt;     &lt;p type="array"&gt;       &lt;_&gt;993330&lt;/_&gt;       &lt;_&gt;993-331&lt;/_&gt;     &lt;/p&gt;   &lt;/contact&gt; &lt;/info&gt; </pre>

Table 3.4: XML to JSON friendly XML(STEP 2)

Algorithm 1	Algorithm 2
<pre> &lt;info type="object"&gt;   &lt;name type="object"&gt;     &lt;f&gt;a&lt;/f&gt;   &lt;/name&gt;   &lt;age&gt;24&lt;/age&gt;   &lt;ismarried&gt;&gt;false&lt;/ismarried&gt;   &lt;city type="object"&gt;     &lt;name&gt;Armonk&lt;/name&gt;   &lt;/city&gt;   &lt;state&gt;NY&lt;/state&gt;   &lt;contact type="object"&gt;     &lt;childtext&gt;home&lt;/childtext&gt;     &lt;p type="array"&gt;       &lt;_&gt;993330&lt;/_&gt;       &lt;_&gt;993-331&lt;/_&gt;     &lt;/p&gt;   &lt;/contact&gt; &lt;/info&gt; </pre>	<pre> &lt;info type="object"&gt;   &lt;name type="object"&gt;     &lt;f type="string"&gt;a&lt;/f&gt;   &lt;/name&gt;   &lt;age type="number"&gt;24&lt;/age&gt;   &lt;ismarried type="boolean"&gt;&gt;false&lt;/ismarried&gt;   &lt;city type="object"&gt;     &lt;name type="string"&gt;Armonk&lt;/name&gt;   &lt;/city&gt;   &lt;state&gt;NY&lt;/state&gt;   &lt;contact type="object"&gt;     &lt;childtext type="string"&gt;home&lt;/childtext&gt;     &lt;p type="array"&gt;       &lt;_ type="number"&gt;993330&lt;/_&gt;       &lt;_ type="string"&gt;993-331&lt;/_&gt;     &lt;/p&gt;   &lt;/contact&gt; &lt;/info&gt; </pre>

**Step 3** After Step 1. and 2., a complete JSON friendly XML is generated. All XML object and array are mapped to  $\langle \text{key}, \text{value} \rangle$  pair of JSON and their respective data types.

```
<info type="object">
  <name type="object">
    <f type="string">a</f>
  </name>
  <age type="number">24</age>
  <ismarried type="boolean">>false</ismarried>
  <city type="object">
    <name type="string">Armonk</name>
  </city>
  <state type="string">NY</state>
  <contact type="object">
    <childtext type="string">home</childtext>
    <p type="array" >
      <_ type="number">993330</_>
      <_ type="string">993-331</_>
    </p>
  </contact>
</info>
```

```
{
  "info":{
    "name":{
      "f":"a"
    },
    "age":24,
    "ismarried":false,
    "city":{
      "name":"Armonk"
    },
    "state":"NY",
    "contact":{
      "childtext":"home",
      "p":[
        993330,
        "993-331"
      ]
    }
  }
}
```

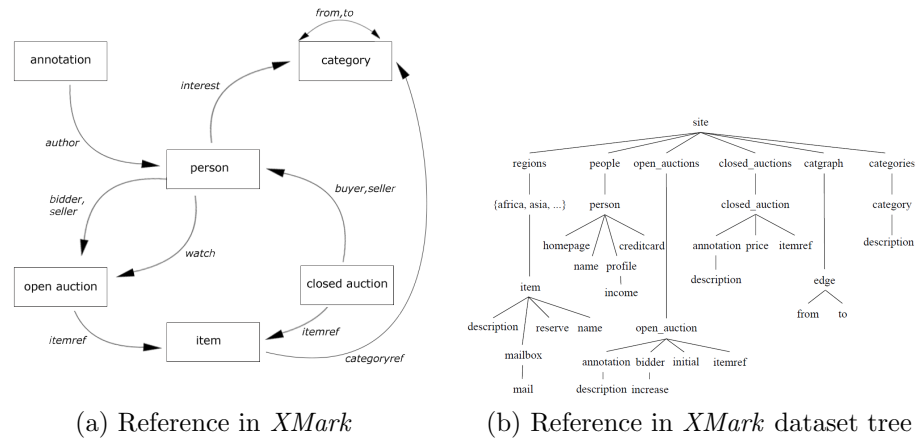
## 3.4 XMark

The XML benchmarking project XMARK (Schmidt et al., 2002) is one of the most popular and most commonly used XML benchmarking project to date. It provides a small executable tool called *xmlgen* that allows for the creation of a synthetic XML dataset based on a fixed schema describing an Internet auctions database. Its generator can be used to build a single record with a large, hierarchical XML tree structure. A factor can be specified to scale the generated data, ranging from a few kilobytes to any arbitrary size limited by the capacity of the system. The textual part of resulting XML document is constructed from 17,000 most frequently occurring words of Shakespeare’s plays.

### 3.4.1 Dataset

The main entities of XMark data come in two groups. *person*, *open\_auction*, *closed\_auction*, *item* and *category* are categorized as first group. The second group entities *annotation* and *description* are natural language text and document-centric element structure. The relationship between the entities in group one are expressed in reference and second group entities are embedded into subtree of group one’s entities. As shown in Fig. 3.1b, XMark dataset has following properties:

- *people* are the collection of *person* that are connected to buyer and seller of open-auctions, closed-auctions, etc. Each person has an unique *id* to reference to another entities like open and closed auctions.
- *items* are the objects for sale or already sold. Each *item* carries a unique identifier and has properties like name, payment information (credit card, money order), description, a reference to the seller, etc; all are encoded as elements. Each item belongs to the world’s regions **africa**, **asia**, **australia**, **europa**, **namerica** and **samerica** as a parent of an item element.
- *open\_auctions* referred to currently in progress auctions which contain bid history(increase/decrease over time) with reference to bidder and seller reference, current bid, the time interval within which the bid is accepted, the status of the transaction, a reference to the item which is being sold etc.
- *closed\_auctions* contains auctions that are successfully completed bidding. They have properties like buyer and seller information, reference to *person*, a reference to an item that is sold, amount of price, quantity of items sold, date of transaction, type of transaction and many more.
- *categories* are used to classify the items. Each category has an unique identifier which is used to reference in item, a name and a description.
- A **catgraph** links categories into a network. The full semantics of XMark dataset can be found in Schmidt et al. (2002).



Q1	Return the name of the person with ID 'person0'.
Q2	Return the initial increase of all open auctions.
Q3	Return the first and current increase of all open auctions whose current increase is at least twice as high as the initial increase.
Q4	List the reserves of those open auctions where a certain person issued a bid before another person.
Q5	How many sold items cost more than 40.
Q6	How many items are listed on all continents?
Q7	How many pieces of prose are in our database?
Q8	List the names of persons and the number of items they bought. (Joins person, closed_auction)
Q9	List the names of persons and the names of items they bought in Europe. (Joins person_auction, item)
Q10	List all persons according to their interest; use French markup in the result.
Q11	For each person, list the number of items currently on sale whose price does not exceed 0.02% of the person's income.
Q12	For each richer-than-average person, list the number of items currently on sale whose price does not exceed 0.02% of the person's income.
Q13	List the names of items registered in Australia along with their description.
Q14	Return the names of all items whose description contains the word 'gold'.
Q15	Print the keywords in emphasis in annotations of closed auctions.
Q16	Return the IDs of those auctions that have one or more keywords in emphasis.
Q17	Which persons don't have a homepage?
Q18	Convert the currency of the reserve of all open auctions to another currency.
Q19	Give an alphabetically ordered list of all items along with their location.
Q20	Group customers by their income and output the cardinality of each group.

Table 3.6: The XMark queries. Source:(Schmidt et al., 2002)



wildcards and includes queries 15 and 16. Category eleven tests the ability to deal with missing values. Only Query 17 is included in this category. The next category deals with user defined functions and contains query 18. Sorting is the purpose category thirteen with single query 19. The last category observe the aggregation with the help of query 20.

### 3.4.3 XMark data into NoSQL Database

A synthetic XMARK dataset consists of one(huge) record in tree structure (Wang et al., 2003). However, as already mentioned in Section 3.4, each subtree in schema, *items*, *people*, *open\_auctions*, *closed\_auctions*, *catgraph* and *categories* contains large number of instances in the database which are indexed. In most of NoSQL system, this scenario is different. Each instance has it's own index structure, therefore the dataset cannot be in just a huge block. As the data model of NoSQL do not match this single structure-encoded sequence, we breakdown it's tree structure into set of sub-structure without losing the overall data and create index for each of them. Beside this, each NoSQL database has their own data model, we need to define model for each of those databases separately.

The generalized concept of XMark data into NoSQL database is given here, but it might be slightly different for each of them. All sub-structures or collections *items*, *people*, *open\_auctions*, *closed\_auctions*, *catgraph* and *categories* are the basic unit for the document fragmentation. Each of these collection stores entities *item*, *person*, *open\_auction*, *closed\_auction* and *category* respectively as mentioned in 3.4.1. These entities represent the documents of NoSQL database. In each of these documents, one special field **type** is added to represent the collection. For example, in case of people substructure, *type* will be *people*. This key value set will be also the part of a document as given in Table 3.8(b). There is one exceptional case for *items* which has **regions** as grandparent and name of different regions like *asia*, *europa* etc. as parent. The *type* for *item* document will be "*items*" as others, one extra field need to be added to represent each region, so there will be one more field *regions* as in Table 3.7.

Table 3.7: Extra added attribute in document

for <i>person</i>	for items which has region name <i>asia</i>
<pre>{   "type": "people" }</pre>	<pre>{   "type": "items",   "regions": "asia" }</pre>

Table 3.8: Example: XMARK data in XML and JSON format

<i>person0</i> in XML(a)	<i>person0</i> in JSON for a NoSQL database(b)
--------------------------	--

<pre> &lt;people&gt;   &lt;person id="person0"&gt;     &lt;name&gt;Kasidit Treweek&lt;/name&gt;     &lt;emailaddress&gt;mailto:Treweek@cohera.com       &lt;/emailaddress&gt;     &lt;phone&gt;+0 (645) 43954155&lt;/phone&gt;     &lt;homepage&gt;http://www.cohera.com/~       Treweek&lt;/homepage&gt;     &lt;creditcard&gt;9941 9701 2489 4716&lt;/       creditcard&gt;     &lt;profile income="20186.59"&gt;       &lt;interest category="category251" /&gt;       &lt;education&gt;Graduate School&lt;/education&gt;       &lt;/profile&gt;     &lt;business&gt;No&lt;/business&gt;   &lt;/person&gt; &lt;/people&gt; </pre>	<pre> {   "id": "person0",   "type": "people",   "name": "Kasidit Treweek",   "emailaddress": "mailto:Treweek@cohera.com",   "phone": "+0 (645) 43954155",   "homepage": "http://www.cohera.com/~Treweek",   "creditcard": "9941 9701 2489 4716",   "profile": {     "income": 20186.59,     "interest": [{       "category": "category251"     }],     "education": "Graduate School",     "business": "No"   } } </pre>
--	---

### 3.4.3.1 XMARK in MongoDB

As mentioned in Section 2.2.1, MongoDB's collections have similar and related documents together. Therefore, for MongoDB, the data modeling concept of 3.4.3 is slightly changed. Each six sub-structure of 3.4.1 are represented as collections. The extra attribute added in previous section *type* is not needed for each document. For *items*, *regions* contains the name of region of each item. Finally, as we mentioned in Indexing, the *id* attribute of each document will be renamed to *\_id* for default indexing. In case of *closed.auctions* and *catgraph*, system will automatically generate *\_id* which is not applicable for us at the moment. A typical example of MongoDB document for person with identifier *person0* is given in Code 3.22.

```

{
  "_id": "person0",
  "name": "Kasidit Treweek",
  "emailaddress": "mailto:Treweek@cohera.com",
  "phone": "+0 (645) 43954155",
  "homepage": "http://www.cohera.com/~Treweek",
  "creditcard": "9941 9701 2489 4716",
  "profile": {
    "income": 20186.59,
    "interest": [{
      "category": "category251"
    }],
    "education": "Graduate School",
    "business": "No"
  }
}

```

Code 3.22: Mongoddb data representation of XMARK data

### 3.4.3.2 Couchbase Server

**Document design** Couchbase Server does not have concept of grouping of documents like *collections* in MongoDB or *tables* of RethinkDB. So, the data model of Section 3.4.3 is applied without alteration. All the documents are stored in a single bucket independently with an *id* as key. For those documents that do not have an *id* attribute, will be manually generated. These documents are identified by *type*. An example of Couchbase Server's document is illustrated in Code 3.23

```

{

```

```

    "id": "person0",
    "type": "people",
    "name": "Kasidit Treweek",
    "emailaddress": "mailto:Treweek@cohera.com",
    "phone": "+0 (645) 43954155",
    "homepage": "http://www.cohera.com/~Treweek",
    "creditcard": "9941 9701 2489 4716",
    "profile": {
        "income": 20186.59,
        "interest": [{
            "category": "category251"
        }],
        "education": "Graduate School",
        "business": "No"
    }
}

```

Code 3.23: Couchbase document of XMARK data for person with id person0

### 3.4.3.3 Rethinkdb

RethinkDB fragment data into tables similar to MongoDB's collection. Each *type* of 3.4.3 is represented as a table. So a XMark database has *items*, *people*, *open\_auctions*, *closed\_auctions*, *catgraph* and *categories* tables each of which content a large numbers of individual document like other NoSQL databases. *id* is the primary index and documents without this field is automatically added an identifier during the time of insertion. In Code 3.24, an illustration of XMark

```

{
    "id": "person0",
    "name": "Kasidit Treweek",
    "emailaddress": "mailto:Treweek@cohera.com",
    "phone": "+0 (645) 43954155",
    "homepage": "http://www.cohera.com/~Treweek",
    "creditcard": "9941 9701 2489 4716",
    "profile": {
        "income": 20186.59,
        "interest": [{
            "category": "category251"
        }],
        "education": "Graduate School",
        "business": "No"
    }
}

```

Code 3.24: Rethinkdb document person0 in people table



# Chapter 4

## Benchmarking

### 4.1 Benchmarking configurations

To produce comparable result it is necessary to run benchmarking test in similar hardware and software configuration. Table 4.1 describe hardware and software information used for Benchmarking.

Table 4.1: Reportable Benchmark Information

Hardware:	CPU: Intel core i3, RAM: 4GB, SSD:250GB
Software:	OS: Linux 3.13.0-49(Ubuntu 14.04), 64 bit
DBMS and Version:	BaseX: 8.1, Mongoddb:2.6.9, Rethinkdb:2.0,Couchbase Server:3.0
Other Software:	Nodejs: v0.12.2, Java:1.7.0_80

### 4.2 Results and Analysis

In this section, We will focus on result of queries on the comparison of different databases. Figure 4.1 represents the time measurements of XMark Queries that are mentioned in Section 3.4.1.

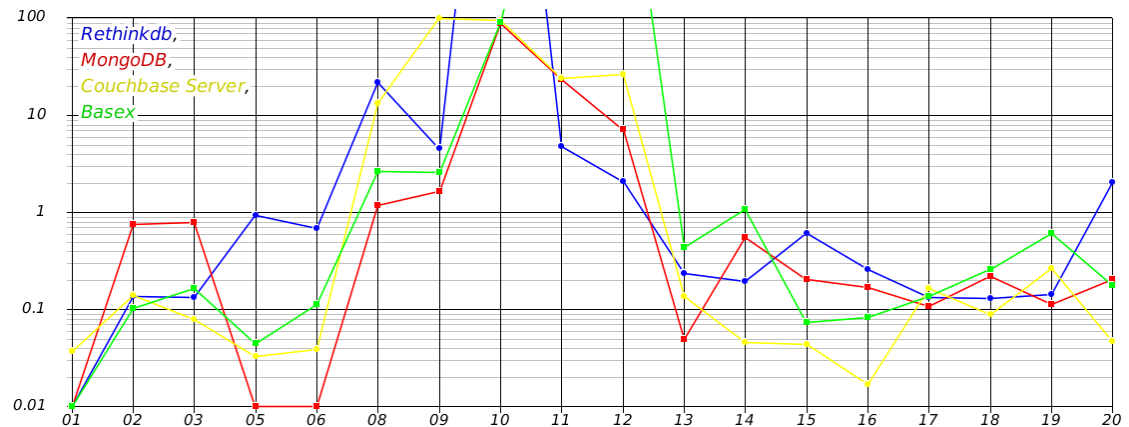


Figure 4.1: Processing of XMark Queries on 111MB XMark instance

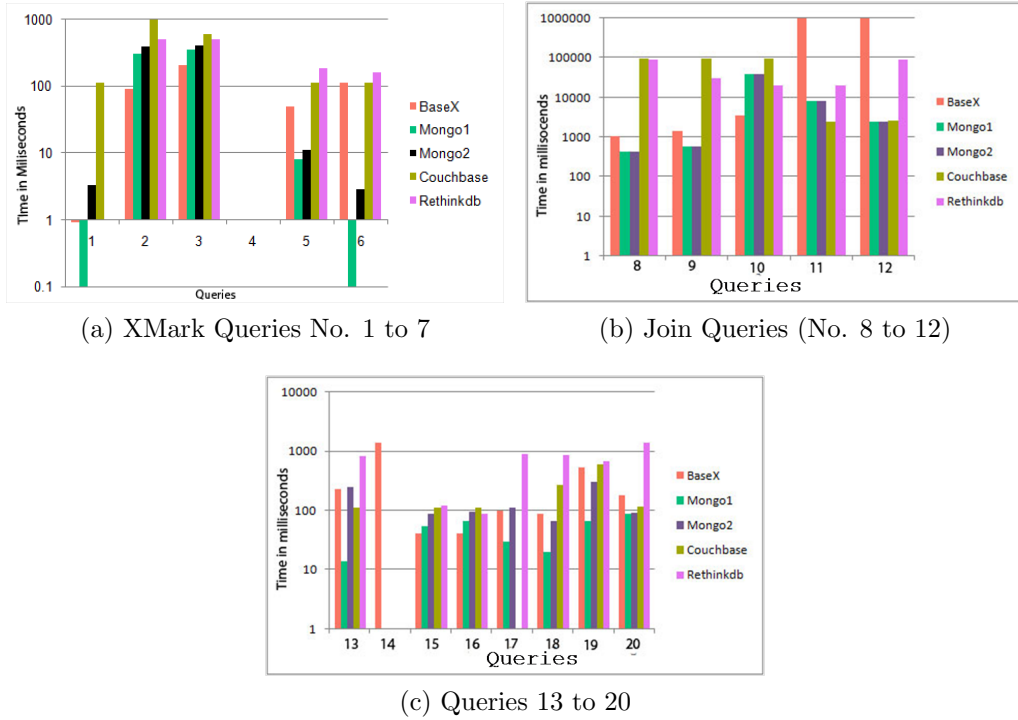


Figure 4.2: XMark Queries time in XML and NoSQL databases(Mongo1: Mongoddb shell, mongo2:through application)

Table 4.2: BaseX

	1	2	3	4	5	6	7	8	9	10
0001	0.0246	0.0201	0.0232	0.0382	0.0303	0.04	0.0377	0.0231	0.0318	0.0252
001	0.0366	0.0374	0.0386	0.0411	0.0315	0.0319	0.0433	0.0989	0.1091	0.1832
01	0.0243	0.0801	0.0946	0.1109	0.0364	0.0347	0.1552	0.4044	0.3664	1.1353
1	0.0352	0.1929	0.1741	0.1281	0.103	0.1876	0.9612	2.3061	2.4278	89.3573
5	0.0181	0.6882	0.7167	0.5554	0.3297	0.5373	4.5321	12.9726	18.1868	2572.0186
10	0.0348	1.6552	1.6249	1.1203	0.7955	1.4824	9.4888	29.3343	37.493	127200
20	0.0318	3.4515	10.1975	3.7436	1.5508	3.6199	27.5222	59.1221	115.4997	12720000.0186
Size	11	12	13	14	15	16	17	18	19	20
0001	0.0142	0.0157	0.0306	0.0365	0.0342	0.0294	0.031	0.036	0.0398	0.0379
001	0.311	0.1782	0.0729	0.0584	0.0335	0.0281	0.0704	0.0361	0.0615	0.0195
01	23.5893	12.5451	0.1074	0.1719	0.0154	0.0144	0.0555	0.0947	0.1527	0.0418
1	2404.1736	1283.4839	0.3643	1.2326	0.0773	0.083	0.1481	0.357	0.5178	0.1857
5	24040.1736	12830.4839	1.6621	6.3384	0.2622	0.203	0.5386	0.6565	2.1216	0.9107
10	240400	128300	2.9848	12.2672	1.2469	0.4025	1.1996	1.9845	5.6113	2.4085
20	240400.1736	128300.4839	5.9743	26.117	2.4337	0.794	2.6421	3.4999	63.6597	7.2226

### 4.2.1 Summary

## 4.3 Discussion

## 4.4 Conclusion

## 4.5 Future Work







# Appendices



# Appendix A

## XMARK Queries

### A.1 MongoDB

Q1. `db.people.find({_id:"person0"},{_id:0,"name":1});`

```
db.open_auctions.aggregate([
  {$project:{_id:1,bidder:"$bidder"}},
  { $unwind: "$bidder"},
  {$group:{_id:"$_id",increase:{$first:"$bidder.
    increase"}}},
  {$sort:{_id:1}},
  { $unwind: "$increase"},
  {$project:{_id:0,increase:1}}
])
```

Q3. `db.open_auctions.aggregate([`  
    `{ $unwind: "$bidder"},`  
    `{$group:{_id:"$_id",first:{$first:"$bidder.`  
        `increase"},last:{$last:"$bidder.increase"}}},`  
    `{$unwind:"$first"},`  
    `{$unwind:"$last"},`  
    `{$project:{_id:1,first:1,last:1,diff:{$gte:["`  
        `$last",{$multiply:["$first",2]}}}}},`  
    `{$match:{diff:true}},`  
    `{$sort:{_id:1}},`  
    `{$project:{_id:0,increase:{first:"$first",last:"`  
        `$last"}}}`  
    `])`

Q4. `[]`

Q5. `db.closed_auctions.find({price:{$gte:40}}).count()`

Q6. `db.regions.find().count()`

Q7. `[]`

Q8.

```

function q8() {
  var ids = new Array();
  var name = new Array();
  var closed = new Array();
  var i=0;
  db.people.find({}, {name:1}).forEach(function(
    people){
    ids[i] = people._id;
    name[ids[i]] = people.name;
    //closed[ids[i]] = {name:people.name,count
      :0};
    i++;
  });
  var closed_auc = db.closed_auctions.aggregate([
    {$project:{_id:1,person:"$buyer.person"}},
    //{$match:{person:{$in:["person12124","
      person12883"]}}},
    {$match:{person:{$in:ids}}},
    {$group:{_id:"$person",count:{$sum:1}}},
    {$sort:{count:-1}},
    {$project:{person:"$_id",count:1,_id:0}}
  ]);
  if(closed_auc) {
    var i=0
    while(closed_auc.hasNext()){
      var auc = closed_auc.next();
      var id = auc.person;
      closed[i] = {"name":name[id],"count
        ":auc.count}
      name[id] = 0;
      i++;
      //printjson(name.indexOf(id));
    }

    }/**/
  var len = closed.length ;
  var j = 0;
  for(var i=0; i < ids.length; i++) {
    if(name[ids[i]] !== 0) {
      //printjson(i);
      closed[len+j] = {"name":name[ids[i]],"
        count":0};
      j++;
    }
  }
  //printjson("Execution Time" + end-start);
  return closed;
}

```

```

Q10.      function q10() {
            var debugId = "person3"
            var ids = new Array();
            var i = 0;
            var allcategories = new Array();
            db.people.aggregate([
                {$match:{"profile.interest":{$exists:true}}},
                {$project: {_id:1, interest:"$profile.interest"}},
                {$unwind:"$interest"},
                {$group: {_id:"$interest.category"}},
                {$project: {_id:0, category:"$_id"}}
            ]).forEach(function(people){
                var catId = people.category;
                ids[i] = {categorie:{id:catId,profile:
                    getProfileByCategory(catId)}}
                i++;
            });
            //printjson("Execution Time" + diff);
            return ids;
        }

```

```

Q11.      function q11() {
            var start = new Date().getTime();
            var debugId = "person3"
            var ids = new Array();
            var open_auc = function(initial){
                return db.open_auctions.find({initial:{$lt:
                    initial}},{_id:1}).count();
            }

            var i=0;
            db.people.find({}, {_id:1, name:1,"profile.income":1}).
                forEach(function(people){
                    var income = ((people.profile) && people.profile.
                        income)? people.profile.income/5000:0;
                    ids[i] = {item:{name: people.name,id:people._id,
                        count:open_auc(income)}};
                    i++;
                });
            return ids;
        }

```

```

Q12.      function q12() {
            var debugId = "person3"
            var ids = new Array();
            var open_auc = function(initial){
                return db.open_auctions.find({initial:{$lt:
                    initial}},{_id:1}).count();
            }
            var i=0;

```

```

db.people.find({"profile.income":{$exists:true}, "profile
.income":{$gt:50000}},{_id:1, name:1,"profile.income
":1}).forEach(function(people){
    var income = ((people.profile) && people.profile.
income)? people.profile.income/5000:0;
    ids[i] = {item:{name: people.name,id:people._id,
count:open_auc(income)}};
    i++;
});
return ids;
}

```

```

    db.regions.find(
        {regions:"australia"},
        {_id:0,name:1,description:1}
    )

```

Q13.

```

    db.regions.aggregate([
        {$match:{$text:{$search:"gold"}}},
        {$group:{$_id:"$_id"}},
        {$group:{$_id:null,count:{$sum:1}}},
        {$project:{$_id:0,count:1}}
    ])

    db.closed_auctions.aggregate([
        {$project:{$_id:1,parlist:"$annotation.description
.parlist"}},
        { $unwind: "$parlist"},
        { $unwind: "$parlist.listitem.parlist"},
        {$project:{$_id:1,text:"$parlist.listitem.parlist.
listitem.text"}},
        {$match:{$or: [{"text.emph.keyword":{$exists:true
}}, {"text.emph.keyword.childtext":{$exists:
true}}, {"text.emph.keyword.child":{$exists:
true}}]}},
        {$project:{$_id:0,text:"$text.emph.keyword"}}
    ])

```

Q15.

```

    db.closed_auctions.aggregate([
        {$project:{$_id:1,parlist:"$annotation.description
.parlist",person:"$seller.person"}},
        { $unwind: "$parlist"},
        { $unwind: "$parlist.listitem.parlist"},
        {$project:{$_id:1,text:"$parlist.listitem.parlist.
listitem.text",person:1}},
        {$match:{"text.emph.keyword":{$exists:true}}},
        {$project:{$_id:0,"person.id":"$person"}}
    ])

```

- Q17. `db.people.find({homepage:{$exists:false}},{_id:0,person  
: "$name"})  
])`
- Q18. `db.open_auctions.aggregate([{$match:{reserve:{$exists:  
true}}},  
{$project:{_id:0,reserve:{$multiply:["$reserve  
",2.20371]}}}])`
- Q19. `db.regions.aggregate([{$sort:{location:1}},  
{$project:{_id:0,"item.name":"$name","item.  
location":"$location"}}])`
- `db.people.aggregate([  
{$project:{_id:1,p:"$profile.income"}},  
{$project:{_id:1,con:  
{$cond:[  
{$gt:["$p",100000]},"preferred",  
{$cond:[{ $and: [{ $lt:["$p", 100000] },  
{$gte:["$p", 30000] } ] },"standard",  
{$cond:[{$and:[{$lt:["$p",30000]},$gt:["$p  
",0]}]},"challenge","na"]}  
}]  
}]  
}},  
{$group:{_id:"$con",sum:{$sum:1}}},  
{$project:{_id:1,sum:1}}])`
- Q20.

## A.2 RethinkDB

- Q1. `r.table("people").get("person0")("name")`
- Q2. `r.table("open_auctions").map(function(doc) {return {  
increase:doc("bidder").nth(0)("increase").nth(0).  
default("")}}})`
- Q3. `r.table("open_auctions").hasFields("bidder").filter(  
r.row("bidder").nth(0)("increase").nth(0).  
mul(2).le(r.row("bidder").nth(r.row("bidder").count().sub(1))("increase").  
nth(0))  
) .map(function(doc){  
var last = doc("bidder").count().sub(1);  
return {increase:{first:doc("bidder").nth  
(0)("increase").nth(0), last:doc("bidder").nth(last)("increase").nth(0)}  
}  
})`

- Q5. 

```
r.table("closed_auctions").count(function(f){
    return f("price").default(0).gt(40)
})
```
- Q6. 

```
r.table("regions").count()
```
- Q8. 

```
r.table("people").map(function(d){
    var x = r.table("closed_auctions").getAll(
        d("id"),{index:"buyer_person"}).count()
    ;
    return {person:d("name"),value:x}
})
```
- Q9. 

```
r.table("regions").getAll("europe",{index:"regions"}).
eqJoin("id",r.table("closed_auctions"),{index:"
closed_itemref_item"}).map({
    name:r.row("left")("name"), id:r.row("
left")("id"),
})
```
- Q10. 

```
r.table("people")("profile")("interest").concatMap(
    function(f){return f}).distinct()
    .merge(function(category){
    return {
        d:
        r.table("people").filter(function(d){
return d("profile")("interest").contains(
    function(e){
return e("category").eq(category("
category"))}))}).map(function(profile
){
    return {
        personne:{
            statistiques:{
                sexe:profile("profile")("gender
").default(""),
                age:profile("profile")("age").
                default(""),
                education:profile("profile")("
education").default(""),
                revenu:profile("profile")("
income").default("")
            },
            coordonnees: {
                nom:profile("name"),
                rue: profile("address")("street
").default(""),
                ville:profile("address")("city
").default(""),
```



```

        pays:profile("address")("
            country").default(""),
        reseau:{
        courrier:profile("emailaddress
            ").default(""),
        pagePerso:profile("homepage").
            default("")
        }
    },
    cartePaiement:profile("creditcard")
        .default("")

    }
}
}).coerceTo('array')
}
})

```

Q11.

```

r.table("people").pluck("id",{profile:{income:true}})
    .map(function(f){
        var i = f("profile")("income").default(0)
            .div(5000);
        return {
            id:f("id"),
            item:r.table("open_auctions").filter(
                function(m){
                    return m("initial").default(0).lt(i)
                }).count()
        }
    })

```

Q12.

```

r.table("people").filter(function(f){return f("profile
    ")("income").default(0).gt(50000)})
    .map(function(f){
        return {
            id:f("id"),
            item:r.table("open_auctions").count(
                function(x){
                    return f("profile")("income").default
                        (0).gt(x("initial").default(0).mul
                            (5000))
                })
        }
    })

```

Q13.

```

r.table("regions").getAll("australia",{index:"regions
    "}).map({item:{name:r.row("name"),description:r.row
        ("description")}})

```

- Q14. 

```
r.table("regions").filter(function(doc) {
    return doc('description').coerceTo("STRING").
        match("gold")
})('name')
```
- Q15. 

```
r.table("closed_auctions")("annotation")("description")
    ("parlist").concatMap(function(f){
    return f('listitem')('parlist')
}).concatMap(function(f){return f('listitem').
    hasFields({text:{emph:{keyword:true}}}))}.map
    (function(f){
    return {"text":f('text')('emph')('keyword')}
    })
```
- Q16. 

```
r.table("closed_auctions").hasFields({annotation:{
    description:{parlist:true}}})
    ("annotation")("description").concatMap(r.
        row("parlist")).hasFields({listitem:{
        parlist:true}})("listitem").concatMap(r.
        row("parlist")).hasFields({listitem:{text
        :{emph:{keyword:true}}}}})
```
- Q17. 

```
r.table("people").filter(function(d){return d.hasFields
    ("homepage").not()}).map({person:{name:r.row("name")
    }})
```
- Q18. 

```
r.table("open_auctions").hasFields("reserve").map(r.row
    ("reserve").mul(2.20371))
```
- Q19. 

```
r.table("regions").orderBy({index:"location"}).pluck("
    location","name").map(function(f){
    return {item:{name:f("name"),location:f("location")}}
    })

r.table("people").map(function(doc){return {income:doc
    ("profile")("income").default(0)}})
    .map({
    preferred:r.branch(r.row("income").ge(100000)
        , 'preferred', 0),
    standard:r.branch(r.row("income").le(100000).
        and(
        r.row("income").ge(30000)
        ), 'standard', 0),
    challenge:r.branch(r.row("income").le(30000).
        and(
        r.row("income").gt(0)
        ), 'challenge', 0),
    na:r.branch(r.row("income").le(0), 'na', 0)
    }).group('preferred', 'standard', 'challenge', 'na')
    .count()
```

## A.3 Couchbase

```
Q1.      function (doc, meta) {
          if(doc.doctype && doc.doctype=="people"
            && doc.id=="person0"){
              emit(doc.name);
          }
      }
```

```
Q2.      function (doc, meta) {
          if(doc.doctype && doc.doctype=="open_auctions"){
              if(doc.bidder){
                  var first = doc.bidder[0];
                  for(var i=0; i<first.length;i++) {
                      if(first[i].increase) {
                          var increase = first[i].increase;
                          //emit(doc.id,increase)
                          emit("increase",increase)
                          break;
                      }
                  }
              } else {
                  emit("increase",null)
              }
          }
      }
```

```
Q3.      function (doc, meta) {
          if(doc.doctype && doc.doctype=="open_auctions"){
              if(doc.bidder){
                  var first = doc.bidder[0];
                  var last = doc.bidder[(doc.bidder.length -
                                          1)];
                  var increaseFirst = 0;
                  var increaseLast = 0;

                  for(var i=0; i<first.length;i++) {
                      if(first[i].increase) {
                          increaseFirst = first[i].increase;
                          //emit(doc.id,increase)
                          break;
                      }
                  }

                  for(var i=0; i<last.length;i++) {
                      if(last[i].increase) {
                          increaseLast = last[i].increase;
                          //emit(doc.id,increase)
                          break;
                      }
                  }
              }
          }
```

```

    }
  }
  if(increaseFirst && increaseLast) {
    if(increaseFirst * 2 <= increaseLast ){
      emit("increase",{first:increaseFirst,last:
        increaseLast,id:doc.id});
    }
  }

  } else {
    // emit("increase",null)
  }
}

```

Q4. []

Q5.

```

function (doc, meta) {
  if(doc.doctype && doc.doctype=="closed_auctions"){
    if(doc.price){
      if(doc.price >= 40) {
        emit(doc.id,doc.price)
      }
    }
  }
}

```

Q6.

```

map:
function (doc, meta) {
  if(doc.doctype && doc.doctype=="regions"){

    emit(doc.id,null)

  }
}

reduce:
_count

```

Q7. []

Q8. []

Q9. []

Q10. []

Q11. []

Q12.           []

Q13.           function (doc, meta) {  
                   if(doc.doctype && doc.doctype == "regions" && doc.  
                       regions && doc.regions=="australia"){  
                       if(doc.description){  
                           emit("item",{name:doc.name, description:doc.  
                               description})  
                       }  
                   }  
           }

Q14.           //don't supprt full text search  
                   function (doc, meta) {  
                       if(doc.doctype && doc.doctype=="regions"){  
                           if(doc.description){  
                               var desc = JSON.stringify(doc.description);  
                               if(desc.indexOf("gold") !== -1)  
                                   emit(meta.id, doc.name);  
                           }  
                       }  
                   }

Q15.           function (doc, meta) {  
                   if(doc.doctype=="closed\_auctions"){  
                       var p1 = doc.annotation.description.parlist;  
                       if(p1){  
                           for(var i=0;i<p1.length;i++){  
                               var p2 = p1[i].listitem.parlist;  
                               if(p2) {  
                                   for(var j=0;j<p2.length;j++){  
                                       if(p2[j].listitem && p2[j].listitem.text &&  
   p2[j].listitem.text.emph && p2[j].listitem  
   .text.emph.keyword){  
   var keyword = p2[j].listitem.text.emph.  
   keyword;  
   emit("text", keyword);  
                                   }  
                               }  
                           }  
                       }  
                   }

Q16.           function (doc, meta) {  
                   if(doc.doctype=="closed\_auctions"){

```

var p1 = doc.annotation.description.parlist;
if(p1){
  var person = (doc.seller.person)?doc.seller.person
    :0;
  for(var i=0;i<p1.length;i++){
    var p2 = p1[i].listitem.parlist;
    if(p2) {
      for(var j=0;j<p2.length;j++){

        if(p2[j].listitem && p2[j].listitem.text &&
          p2[j].listitem.text.emph && p2[j].listitem
            .text.emph.keyword){
          var keyword = p2[j].listitem.text.emph.
            keyword;
          if(person){
            emit("person", {id:person});
          }
        }
      }
    }
  }
}

```

Q17.

```

function (doc, meta) {
  if(doc.doctype=="people"){
    if(!doc.homepage && doc.name){
      emit("name",doc.name)
    }
  }
}

```

Q18.

```

function (doc, meta) {
  if(doc.doctype=="open_auctions"){

    if(doc.reserve){
      var exchange = parseFloat(doc.reserve) * 2.20371;

      emit(null,exchange)
    }
  }
}

```

Q19.

```

function (doc, meta) {
  if(doc.doctype=="regions" && doc.name){
    var location = (doc.location)?doc.location:"";
    emit(location,{item:{name:doc.name,location:
      location}});
  }
}

```

```
    }
  }
})

map:
function (doc, meta) {
  if(doc.doctype=="people"){

    var income = (doc.profile && doc.profile.income)?
      doc.profile.income : 0;

    if(income >= 100000 ){
      emit("preferred",1);
    }else if(income < 100000 &&
      income >= 30000) {
      emit("standard",1);
    }else if(income < 30000 &&
      income > 0 ){

      emit("challenge",1);
    } else {
      emit("na",income);
    }
  }
}

reduce:
_count
```





# Bibliography

- Q20Veronika Abramova and Jorge Bernardino. Nosql databases: Mongodb vs cassandra. In *Proceedings of the International C\* Conference on Computer Science and Software Engineering*, pages 14–22. ACM, 2013.
- Rick Cattell. Scalable sql and nosql data stores. *ACM SIGMOD Record*, 39(4): 12–27, 2011.
- Xiaoming Gao. Investigation and comparison of distributed nosql database systems. 2013.
- Christian Grün. BaseX – The XML Database for Processing, Querying and Visualizing large XML data. <http://basex.org>, October 2010.
- Christoforos Hadjigeorgiou et al. RDBMS vs NoSQL: Performance and scaling comparison. *The University of Edinburgh*, 2013.
- Robin Hecht and S Jablonski. Nosql evaluation. In *International Conference on Cloud and Service Computing*, 2011.
- Haifeng Jiang, Hongjun Lu, Wei Wang, and Jeffrey Xu Yu. XParent: An efficient RDBMS-based XML database system. In *Data Engineering, 2002. Proceedings. 18th International Conference on*, pages 335–336. IEEE, 2002.
- David Lee. Jxon: an architecture for schema and annotation driven json/xml bidirectional transformations. In *Proceedings of Balisage: The Markup Conference*, 2011.
- Kyle Lichtenberg and Miles Ward. Nosql database in the cloud: Couchbase server 2.0 on aws. 2013.
- Irena Mlýnková. Xml benchmarking: Limitations and opportunities. Technical report, Technical Report, Department of Software Engineering, Charles University, Czech Republic, 2008.
- Cory Nance, Travis Losser, Reenu Iype, and Gary Harmon. NoSQL vs RDBMS-why there is room for both. In *Proceedings of the Southern Association for Information Systems Conference*, pages 111–116, 2013.
- Kai Orend. Analysis and classification of nosql databases and evaluation of their ability to replace an object-relational persistence layer. *Architecture. Citeseer*, page 100, 2010.
- David Ostrovsky and Yaniv Rodenski. *Pro Couchbase Server*. Apress, 2014.
- Rethinkdb. *An open-source distributed database built with love - RethinkDB*, 2015 (accessed January 12, 2015). URL <http://www.rethinkdb.com/>.

- Albrecht Schmidt, Florian Waas, Martin Kersten, Michael J Carey, Ioana Manolescu, and Ralph Busse. Xmark: A benchmark for xml data management. In *Proceedings of the 28th international conference on Very Large Data Bases*, pages 974–985. VLDB Endowment, 2002.
- Jayavel Shanmugasundaram, Kristin Tufte, Chun Zhang, Gang He, David J DeWitt, and Jeffrey F Naughton. Relational databases for querying XML documents: Limitations and opportunities. In *Proceedings of the 25th International Conference on Very Large Data Bases*, pages 302–314. Morgan Kaufmann Publishers Inc., 1999.
- Haixun Wang, Sanghyun Park, Wei Fan, and Philip S. Yu. Vist: A dynamic index method for querying xml data by tree structures. In *In SIGMOD*, pages 110–121, 2003.
- Cong Yu and HV Jagadish. Schema summarization. In *Proceedings of the 32nd international conference on Very large data bases*, pages 319–330. VLDB Endowment, 2006.

# List of Figures

2.1	column based database . . . . .	4
2.2	A document in document based database(JSON) . . . . .	4
2.3	MongoDB document structure . . . . .	5
2.4	MongoDB Clusters . . . . .	8
2.5	Couchbase Server vBucket . . . . .	9
2.6	Couchbase Server's Document Design . . . . .	10
3.1	XMark data tree and reference (Schmidt et al., 2002) . . . . .	19
3.2	XMark ER-Diagram. Nodes, solid arrows, and dashed arrows represent schema elements (or attributes, with prefix '@'), structural links, and value links, respectively. Elements with suffix '*' are of SetOf type(Yu and Jagadish, 2006) . . . . .	19
4.1	Processing of XMark Queries on 111MB XMark instance . . . . .	25
4.2	XMark Queries time in XML and NoSQL databases(Mongo1: MongoDB shell, mongo2:through application) . . . . .	26



# List of Tables

2.1	Mapreduce in Couchbase Server . . . . .	10
3.1	Anonymous values of JSON in XML . . . . .	13
3.2	Translation of simple XML data types into JSON . . . . .	14
3.3	XML to JSON friendly XML(STEP 1) . . . . .	16
3.4	XML to JSON friendly XML(STEP 2) . . . . .	16
3.5	XML to JSON (STEP 3.) . . . . .	16
3.6	The XMark queries. Source:(Schmidt et al., 2002) . . . . .	20
3.7	Extra added attribute in document . . . . .	21
3.8	Example: XMARK data in XML and JSON format . . . . .	21
4.1	Reportable Benchmark Information . . . . .	25
4.2	BaseX . . . . .	26
4.3	MongoDB With Index . . . . .	27
4.4	MongoDB Without Index . . . . .	27
4.5	RethinDB . . . . .	27
4.6	Couchbase . . . . .	27



# Listings

2.1	MongoDB sample document . . . . .	6
2.2	. . . . .	10
2.3	. . . . .	10
2.4	Chainable Query in ReQL . . . . .	11
3.1	. . . . .	13
3.2	. . . . .	13
3.3	. . . . .	13
3.4	. . . . .	14
3.5	. . . . .	14
3.6	. . . . .	14
3.7	. . . . .	14
3.8	. . . . .	14
3.9	. . . . .	14
3.10	. . . . .	14
3.11	. . . . .	14
3.12	. . . . .	16
3.13	. . . . .	16
3.14	. . . . .	16
3.15	. . . . .	16
3.16	. . . . .	17
3.17	. . . . .	17
3.18	. . . . .	21
3.19	. . . . .	21
3.20	. . . . .	22
3.21	. . . . .	22
3.22	Mongodb data representation of XMArk data . . . . .	22
3.23	Couchbase document of XMArk data for person with id person0 . . . . .	22
3.24	Rethinkdb document person0 in people table . . . . .	23