



From Relational to Riak

December 2012

Table of Contents

| | |
|---|----------|
| Table of Contents | 1 |
| Introduction | 1 |
| Why Migrate to Riak? | 1 |
| The Requirement of High Availability | 1 |
| Minimizing the Cost of Scale | 2 |
| Simple Data Models | 4 |
| Tradeoff Decisions | 5 |
| Eventual Consistency | 5 |
| Data Modeling | 5 |
| Operational & Development Considerations | 6 |
| Data Migration | 6 |
| The Basics of Modeling Data In Riak | 7 |
| Resolving Data Conflicts | 8 |
| Multi-Datacenter Operations | 8 |
| Conclusion | 9 |

Introduction

This technical brief is designed to provide a background and detailed level of understanding for those analyzing a move from a relational database to a NoSQL model (specifically emphasizing the Riak key/value store offered by Basho).

The brief begins by examining commonly cited reasons for choosing Riak instead of a relational database, specifically focused on availability vs. consistency tradeoffs, scalability, and the key/value data model. Then it analyzes the decision points that should be considered when choosing a non-relational solution and what such offerings cannot provide related to querying, data modeling, and consistency guarantees. Finally, it provides simple patterns for building common applications in Riak using its key/value design; dealing with data conflicts that emerge in an eventually consistent system; and how replicating data to multiple sites is possible with Riak.

Why Migrate to Riak?

This section analyzes the most common reasons to move from a relational database to Riak.

The Requirement of High Availability

Relational Databases tend to favor consistency over availability, making them ill suited for applications that require high availability

The CAP theorem, fathered by Dr. Eric Brewer, states that in the event of a network partition, a distributed system can either provide availability or consistency. Relational systems typically leverage master/slave architecture in order to ensure consistency of operations – that all replicas receive the same data at the same time. Consistent operations provide applications with guarantees that read operations reflect the last successful update to the database, and are an important facet of enabling operations, like transactions, that are essential for some types of applications (billing and financial systems being the canonical examples).

Relational databases are, most commonly, seen in production running on a single server; however, if the dataset grows beyond the capacity of a single machine it becomes necessary to scale the database and operate in a distributed environment. Relational databases address the challenge of scale with a master/slave architecture, wherein the topology of a cluster is comprised of a single master node and multiple slaves. This architecture utilizes a technique called “sharding” to achieve scale. Sharding breaks data into parts that can be distributed across multiple physical machines. A common example would be putting user data for differing geographical regions (e.g., US and EU) on different machines, or using an alphabetical or numerical order to split data.

Under this configuration, the master node is responsible for accepting all write operations and coordinating with slave nodes to apply the updates in a consistent manner. Read requests can either be proxied through the master or sent directly to the responsible slave.

However, in the event that a master node fails, the database will favor consistency and reject write operations until the failure is resolved. This can lead to a window of write unavailability, which is unacceptable in some application designs. Most master/slave architectures recognize that a master node is a single point of failure and will perform automatic master failover, wherein a slave will be elected as a new master when failure of the master node is detected.

In contrast, Riak is a masterless system designed to favor availability even in the event of node failures and/or network partitions. Any node can serve any incoming request, regardless of data locality. Each object in Riak is replicated a configurable number of times in the cluster (the default is three). If a node experiences an outage, other nodes will continue to server write and read request. Further, if a node becomes unavailable to the rest of the cluster, a neighboring node will take over write and update responsibilities for the “missing” node. The neighboring node will pass new, or update, objects back to the original node once it rejoins the cluster through a process called “hinted handoff”.

Riak’s masterless design ensures read and write availability is maintained even if many nodes become unavailable due to network partition or hardware failure. However, a lack of master/slave configuration to ensure consistency means that data in Riak is *eventually consistent* – all updates will eventually propagate to all nodes.

For many of today’s application and platforms, high availability is more important than strict consistency. In some use cases, Data unavailability can have a direct impact on revenue – cloud services, online retail, shopping carts, checkout process, advertising are just a few examples. Further, lack of availability can damage user trust and result in a poor use experience for many websites, social, and mobile applications; and for critical data, like user data or serving content, availability can be more important than strict consistency.

Minimizing the Cost of Scale

Scaling a relational database to handle more data and usage can be prohibitively expensive for operators.

As described above, relational databases generally leverage a technique called “sharding” to achieve scale requirements. Sharding breaks data into parts that can be distributed across multiple physical machines.

This approach can be problematic for several reasons. First, writing and maintaining sharding logic increases the overhead of operating and developing an application on the database. Significant growth of data or traffic typically means significant, often manual, resharding projects. Determining how to intelligently split the dataset without negatively impacting performance, operations, and development presents a substantial challenge, especially when dealing with “big data”, rapid scale, or peak loads. Further, rapidly growing applications frequently outpace an existing sharding scheme. When the data in a shard grows too large, the shard must again be split. While several “auto”-sharding technologies have emerged in recent years, these methods are often imprecise and manual intervention is standard practice. Finally, sharding can often lead to “hot spots” in the database – physical machines with a disproportionately high amount of data or serving a disproportionately high number of request – which can lead to uneven latency and degraded performance.

In Riak, data is distributed evenly across nodes using consistent hashing. Consistent hashing ensures data is event distributed around the cluster and new nodes can be added with automatic, minimal reshuffling of data. This significantly decreases risky “hot spots” in the database and lowers the operational burden of scaling.

Riak stores data using a simple key/value scheme. These keys and values are stored in a namespace called a bucket. When new key/value pairs are added to a bucket, object’s bucket/key pair is hashed. The resulting value maps onto a 160-bit integer space. This integer space can be conceptualized as a ring that is used to determine where data is placed on physical machines.

Riak tokenizes the total key space into a fixed number of equally sized partitions (default is 64). Each partition owns the given range of values on the ring and is responsible for all buckets and keys that, when hashed, fall into that range. A virtual node (or “vnode”) is the process that manages each of these partitions. Physical machines evenly divide responsibility for vnodes.

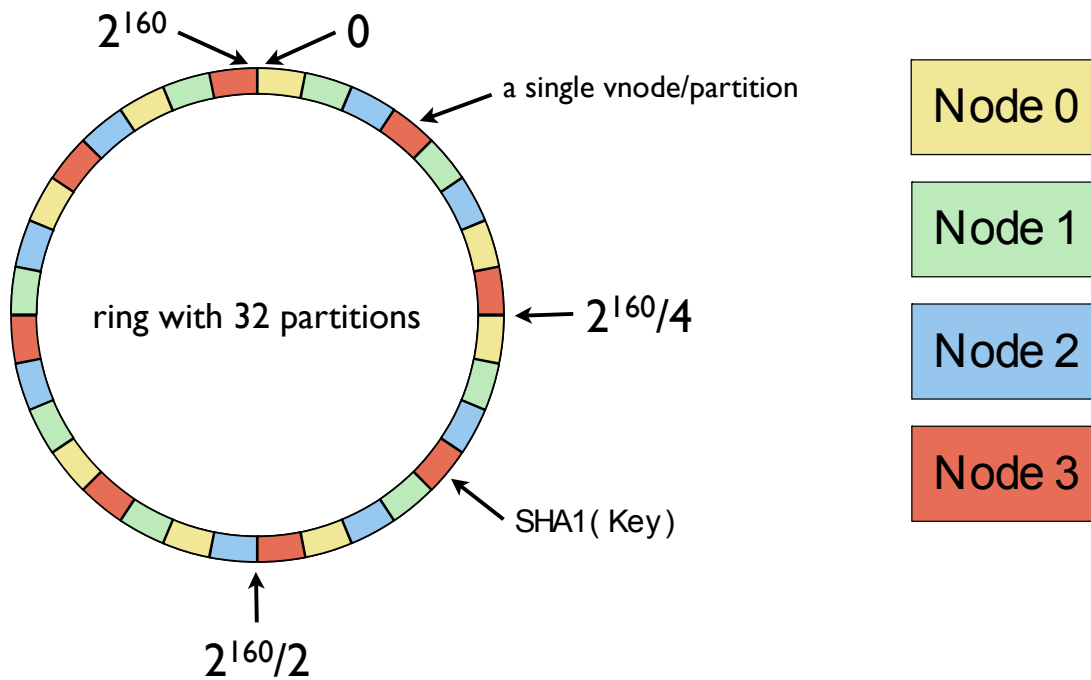


Figure 1: The Riak "Ring"

As a result of the even distribution created by the hashing function, and how physical nodes share responsibility for keys, Riak ensures data is evenly dispersed. When machines are added, data is rebalanced automatically with no downtime. New machines take responsibility for their share of data by assuming ownership of some of the partitions; existing cluster members hand off the relevant partitions and the associated data. The new node continues claiming partitions until data ownership is equal, updating a picture of which nodes own what as it goes. This picture of cluster state is shared to every node using a gossip protocol and serves as a guide to route requests. This process is what ensures that any node in the cluster is able to receive requests, ensuring that developers do not need to worry about the underlying complexity of identifying where the data lives.

For many Riak users, consistent hashing, and the resulting model for adding capacity as needed, provides a much simpler operational scenario than manually sharding data.

Simple Data Models

The relational data model can be needlessly complex and inflexible for certain types of applications.

The data models of traditional Relational Database Management Systems (RDBMS) offer many features that developers find important. However, with the emergence of trends like big data, “agile” development, and new types of applications (social & mobile), there has also been a significant growth of unstructured data, data that does not require the rigid data model of relation systems.

Many applications can effectively utilize a simple key/value model for storing and retrieving data. In Riak, objects are comprised of key/value pairs, which are stored in flat namespaces called “buckets”. Riak does not dictate what types of data are persisted – all objects are stored on disk as binaries. As a result, developing code that interacts with this simple, straightforward design can be a faster process. Additionally, adding new features to the application will not require updating a scheme or changing the data model, ideal for applications where rapid iterations is required and changes in the underlying data store are undesirable. A key/value design is not appropriate for all applications and use cases, but for many this model provides more flexibility and simplicity thereby contributing to developer productivity. Later in this technical brief, some common use cases for Riak, and approaches to modeling data for those use cases, are discussed. This includes content, advertising, sensor, user, and log data.

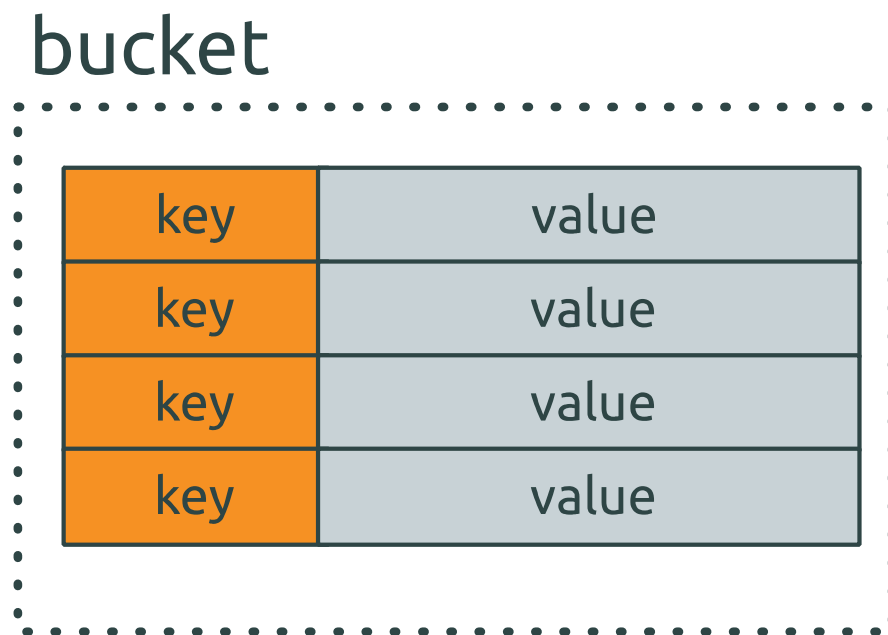


Figure 2: Key/Value Pairs Stored in a Bucket

Tradeoff Decisions

Eventual Consistency

Riak does not support consistent operations

Data that must be consistent – because the application needs to perform transactions, requires a strict guarantee around ordering, or other reasons – is not a good fit for Riak.

Riak does provide some tunable controls for correctness vs. availability – on a per-request basis, users can set an w and r value that adjust the number of replicas that must respond to a read or write request for it to succeed. For example, if Riak has been configured to store three replicas of each object, a request with a w value of 3 requires all three replicas to acknowledge the write before it is considered successful. A “durable write” value, also on a per-request basis, indicates how many replicas must commit to durable storage before returning a success result. However, despite these configurable settings, Riak cannot guarantee strict consistency (especially in the presence of failure conditions) and should not be used for data that require such strict guarantees.

Data Modeling

Riak’s design does not support the richer data types and model of traditional relational systems

While many users find that Riak’s key/value model is more flexible, faster to develop against, and well suited to their applications, there are tradeoffs regarding query options and data types available. Riak does not expose sets, counters, or transactions; it does not support join operations and there is no concept of columns and rows. Riak is queried via HTTP request, via the protocol buffers API, or through various client libraries; there is no SQL or SQL-like language. Riak’s simpler data model results in fewer and leaner query ability options.

Riak does, however, offer additional functionality on top of the fundamental key/value model:

- **Riak Search:** Riak Search is a distributed, full-text search engine. It provides a SOLR-like API, support for various MIME types & analyzers, and robust querying including exact matches, wildcards, range queries, proximity searches, and more.
- **Secondary Indexing:** Secondary Indexing (2i) in Riak gives developers the ability, at write time, to tag an object stored in Riak with one or more queryable values. Indexes can be either integers or strings, and can be queried by either exact matches or ranges of an index.
- **MapReduce:** Developers can leverage Riak MapReduce for tasks like filtering documents by tag, counting words in documents, and extracting links to related data. It offers support for both JavaScript and Erlang.

For more information, check out the Riak documentation on [Querying Data](#). Additionally, work is being done to expose additional data types that are tolerant of Riak's eventually-consistent nature, specifically counters and sets, however this is not yet available in a public release.

Operational & Development Considerations

Data Migration

How should data be migrated to Riak? The topic of data migration could fill an entire book, so these are just a few points of starting advice. Should you desire in-depth help or consultation, the [Professional Services](#) team at Basho is always available for assistance.

The recommended method of migration is to adopt a staged approach, migrating areas of the application to Riak while running it alongside the previous data storage mechanism. For each stage, pick a standalone logical unit of data (a group of related tables, or a document in the current storage system), convert it to a storage format appropriate to Riak (the columns of a relational table map easily into JSON or XML fields), consider how the data will be accessed (will it be plain key/value reads and writes, or will it involve more complex queries), and write migration scripts. The actual movement of the data will involve code in the preferred programming language. It is necessary to iterate over existing data (either by connecting directly via a database driver, or reading in export CSV or XML data from disk), compose appropriate objects for storage, and use Riak client code to store those objects in the Riak cluster.

Start with areas of data that have one-to-one relationships, direct lookup on keys that are already known. Most applications will have standalone high-traffic areas that are perfect to model as key/value storage operations, such as sessions, user preferences, advertisements, small binary or XML/JSON documents. A sure sign of these areas is that the applications gets the keys "for free", without having to perform any queries to discover them. For example, a user logs in and gets assigned a session cookie – the application now has two easy keys in memory (a user id and session id) with which to load data. Or, a web request comes in and returns an id as part of the URL, such as a request for a content node in a wiki or CMS, or a request to serve a particular ad.

If a relation database is already in use, these will be that objects that are loaded from a single table, or from several tightly related ones. They will, likely, have already been optimized for high read traffic, possibly de-normalized from several related tables, for ease and speed of access. If possible, do not choose an area that requires atomic multi-step transactions the first point of migration; the topic of consistency and collision handling in highly parallel distributed systems is an advance topic that requires a deep understanding of the workings of Riak and the capabilities it provides.

Once a slice of the data model has been isolated for migration, consider how it will be stored in Riak – what key and object format will be used. In most cases, the keys will be dictated by the existing application data (the format of the session id or user id will be fixed), and these objects can be reused as Riak object keys. The format of the object payload requires additional consideration. If small binaries are being stored (PDF documents, small images, or custom binary data objects), these can be stored directly as binary blobs. If structured data is being migrated (for example, relational database tables), consider using structured text documents such as JSON or XML. If metadata is required alongside the object (timestamps, owner ids, tags of any kind), consider its storage location – as custom Riak object headers, or as additional fields in the JSON/XML object payload.

Having made these design choices, the migration becomes fairly straightforward. Code must be written in a preferred programming language to obtain the data from an existing system, composed into appropriate Riak objects with the chosen keys and values, and writes issued to the Riak cluster. There is no automated or “backend” way to migrate the data into Riak; the only way to ensure data is properly stored, and the appropriate indexes are created (in the case of leverage Search or Secondary Index functionality), is to use the Riak client API to perform the writes.

After migration of the obvious areas perfect for simple key/value reads and writes, it is likely desired to migrate the data that requires more advanced querying capability or complex one-to-many and many-to-many relationships. [Basho documentation](#) and tutorials provide additional information on data modeling and use cases, Link Walking, Key Filtering, Secondary Indexes, Search, and Map/Reduce queries. While there is an initial learning curve when transitioning from a traditional relational model to a distributed key/value system, and there is no exact equivalent to the familiar table join, rest assured that Riak does offer powerful and flexible data modeling and querying capabilities that have been field-tested in product projects of massive scale. [The Riak-users mailing list](#) is an available source of information, or [contact Basho directly](#) for advice, code and architecture review, and consultation.

The Basics of Modeling Data In Riak

Below is a chart with some simple approaches to building common application types with a key/value model. Remember that values in Riak are opaque and stored on disk as binaries – JSON or XML documents, images, text, etc. It is worth noting that that how an application is structured to run on Riak should take into account the unique needs of the application, includes access patterns such as read/write distribution, latency differences between various operations, use of Riak features including MapReduce, Search, Secondary Indexes, and more.

| Application Type | Key | Value |
|--------------------|--------------------|--|
| Session | User/Session ID | Session Data |
| Advertising | Campaign ID | Ad Content |
| Logs | Date | Log File |
| Sensor | Date, Date/Time | Sensor Updates |
| User Data | Login, eMail, UUID | User Attributes |
| Content | Title, Integer | Text, JSON/XML/HTML Document, Images, etc. |

For additional information, and more complex considerations such as modeling relationship and advance social applications, see the Riak documentation on [use cases and data modeling](#).

Resolving Data Conflicts

In any system that replicates data, conflicts can arise – e.g., if two clients update the same object at the exact same time; or if not all updates have yet reached hardware that is experiencing lag. As discussed earlier, Riak is “eventually consistent” – while data is always available, not all replicas may have the most recent update at the exact same time, causing brief periods (generally on the order of milliseconds) of inconsistency while all state changes are synchronized.

However, Riak does provide features to detect and help resolve the statistically small number of incidents when data conflicts occur. When a read request is performed, Riak looks up all replicas for that object. By default, Riak will return the most updated version, determined by looking at the object’s vector clock. Vector clocks are metadata attached to each replica when it is created. They are extended each time a replicate is updated to keep track of versions. Clients can also be allowed to resolve conflicts themselves.

Further, when an outdated object is discovered as part of a read request, Riak will automatically update the out-of-sync replica to make it consistent. Read repair, a self-healing property of the database, will even update a replica that returns a “not_found” in the event that a node loses it due to physical failure.

For information on vector clocks and conflict resolution can be found in the [online documentation](#).

Multi-Datacenter Operations

Multi-site replication is quickly becoming critical for many of today’s platform and applications. Not only does replication across multiple clusters provide geographic data locality – the ability to serve global traffic at low-latencies, it can also be an integral part of a disaster recovery or backup strategy. Other teams may use multi-site replication to maintain secondary data stores for both failover as well as for performing intensive computation without disrupting production load. Multi-site replication is included in Basho’s commercial extension to Riak, [Riak Enterprise](#), which includes 24/7 support.

Multi-site replication in Riak works differently than the typical approach seen in the relational world, specifically multi-master replication. In Riak’s multi-datacenter replication, one cluster acts as a “primary cluster”. The primary cluster handles replication request from one or more “secondary clusters” (generally located in datacenter in other regions or countries). If the datacenter with the primary cluster goes down, a secondary cluster can take over as the primary cluster. In this sense, Riak’s multi-datacenter capabilities are “masterless”.

In multi-datacenter replication, there are two primary modes of operation: full-sync and real-time. In full-sync mode, a complete synchronization occurs between primary and secondary cluster(s). In real-time mode, continual, incremental synchronization occurs – replication is triggered by new updates. Full-sync is performed upon initial connection of a secondary cluster, and then periodically (by default, every 360 minutes). Full-sync is also triggered if the TCP connection between primary and secondary clusters is severed and then recovered.

Data transfer is unidirectional (primary->secondary). However, bidirectional synchronization can be achieved by configuring a pair of connections between clusters.

Full documentation for multi-datacenter replication in Riak Enterprise is available in the [online documentation](#).

Conclusion

Picking the right database for your team means a careful understanding of the requirements of your application or platform, what developer productivity means to you, the operational conditions you need, and how different database solutions support those goals. We hope this gets you started with understanding the differences between traditional database solutions and Riak, and how you can be successful in a non-relational world.

If you are interested in hearing more about Riak users and use cases, review the [documented user stories](#).

If you have additional questions, [get in touch](#) – the Basho team would be happy to discuss your use case and help determine if Riak is the appropriate technological fit.