



POLITECNICO  
MILANO 1863

SYSTEMS AND METHODS FOR BIG AND UNSTRUCTURED DATA

# Key-value Databases – Redis

Marco Brambilla

[marco.brambilla@polimi.it](mailto:marco.brambilla@polimi.it)

 @marcobrambi

# Overview – NoSQL Family

Data stored in 4 types:

- Document
- Graph
- Key-value
- Wide-column

Document Database	Graph Databases
 	 <b>Neo4j</b>  <b>InfiniteGraph</b> The Distributed Graph Database
Key-value Databases	Wide Column Stores
 <b>redis</b>  <b>AEROSPIKE</b>  <b>amazon DynamoDB</b>  <b>riak</b>	 <b>ACCUMULO</b>  <b>HYPERTABLE</b>  <b>APACHE HBASE</b>  <b>Amazon SimpleDB</b>

# Why is performance important?

- Amazon - Every 1/10 second delay resulted in 1% loss of sales.
- Google - Half a second delay caused a 20% drop in traffic.
- Industrial Group - 1-second delay in page-load time
  - 11% fewer page views
  - 16% decrease in customer satisfaction
  - 7% loss in conversions

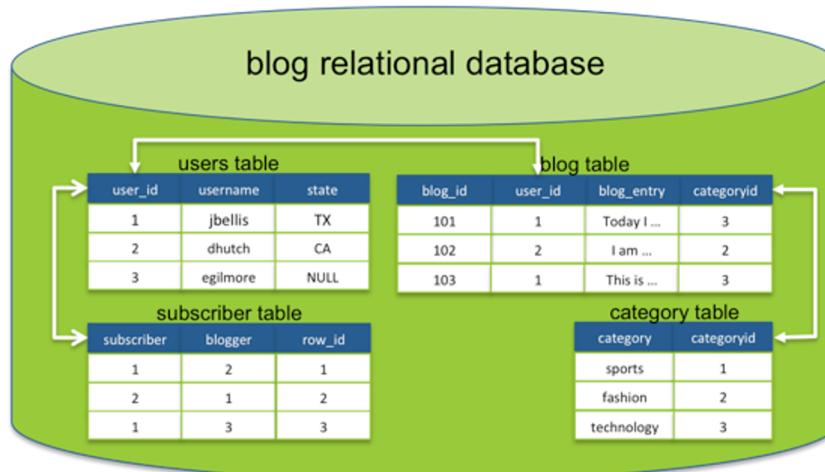
# Key-Value databases

# Why Key-value Store?

- (Business) Key -> Value
- (twitter.com) tweet id -> information about tweet
- (kayak.com) Flight number -> information about flight
- (yourbank.com) Account number -> information about it
- (amazon.com) item number -> information about it
- Search by ID is usually built on top of a key-value store

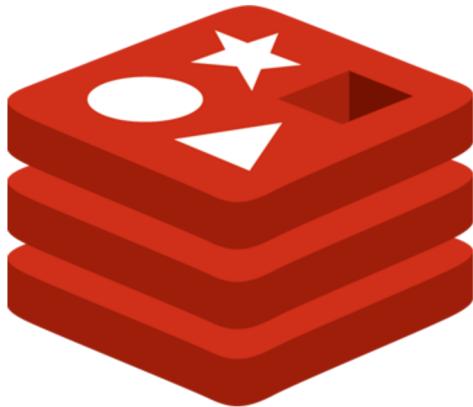
# Isn't that just a database?

- Queried using SQL
- Key-based
- Foreign keys
- Indexes
- Joins



SQL queries: `SELECT user_id from users WHERE username = "jbellis"`

[Example's Source](#)



redis  
REmote DIctionary Server

# What is REDIS ?

- ▶ Redis is an advanced **key-value store**, where keys can contain data structures such as **strings**, **hashes**, **lists**, **sets**, and **sorted sets**. Supporting a set of **atomic operations** on these data types.
- ▶ Redis is a different evolution path in the key-value databases where values are complex data types that are closely related to fundamental data structures and are exposed to the programmer as such, without additional abstraction layers.
- ▶ Can be used as **Database**, a **Caching layer** or a **Message broker**.

Redis can persist data to the disk

Redis is fast

Redis is not only a key-value store

# What is not REDIS

- ▶ Redis is not a replacement for Relational Databases nor Document Stores.
- ▶ It might be used complementary to a SQL relational store, and/or NoSQL document store.
- ▶ Even when Redis offers configurable mechanisms for persistency, increased persistency will tend to increase latency and decrease throughput.
- ▶ *Best used for rapidly changing data with a foreseeable database size (should fit mostly in memory).*

*NoSQL comparisons:*

<http://kovacs.eu/cassandra-vs-mongodb-vs-couchdb-vs-redis>

<http://www.infoivy.com/2013/07/nosql-database-comparison-chart-only.html>

# Redis Use Cases

- ▶ Caching
- ▶ Counting things
- ▶ Blocking queues
- ▶ Pub/Sub (service bus)
- ▶ MVC Output Cache provider
- ▶ Backplane for SignalR
- ▶ ASP.NET Session State provider\*
- ▶ Online user data (shopping cart, ...)
- ... Any real-time, cross-platform, cross-application communication

\* ASP.NET session state providers comparison: <http://www.slideshare.net/devopsguys/best-performing-aspnet-session-state->

# When to consider Redis

- Speed is critical
- More than just key-value pairs
- Dataset can fit in memory
- Dataset is not critical

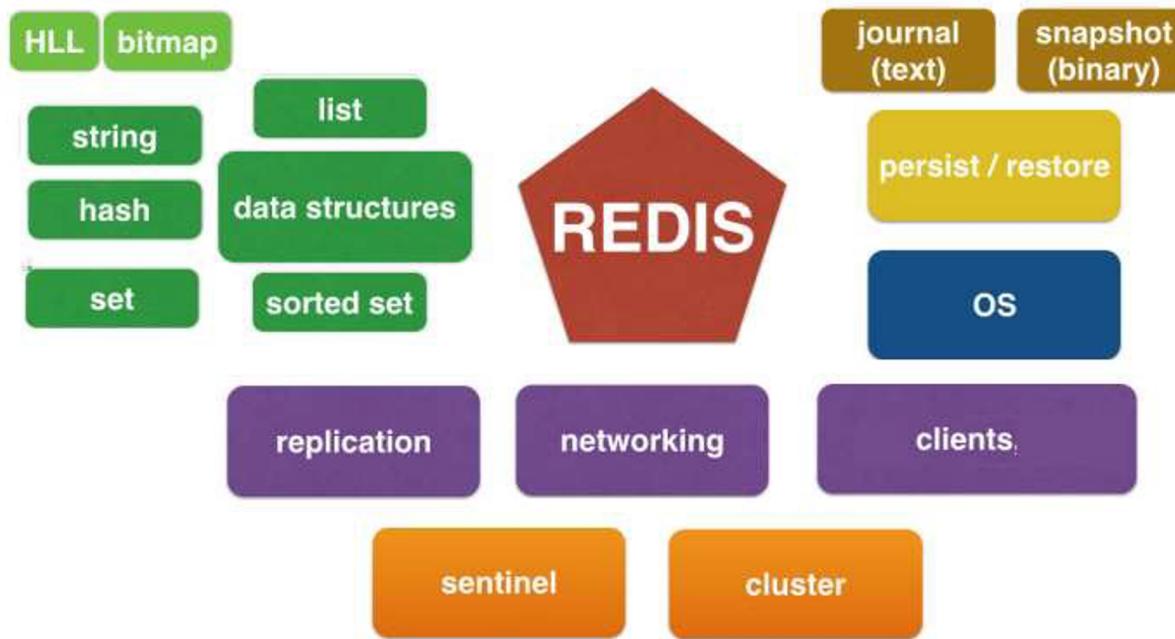
# Redis Architecture

- ▶ Written in ANSI C by Salvatore Sanfilippo (@antirez).
- ▶ Works in most POSIX systems like Linux, BSD and OS X.
- ▶ Linux is the recommended
- ▶ No official support for Windows, but Microsoft develops and maintains an open source Win-64 port of Redis\*
- ▶ Redis is a single-threaded server, not designed to benefit from multiple CPU cores.  
Several Redis instances can be launched to scale out on several cores.
- ▶ All operations are atomic (no two commands can run at the same time).
- ▶ It executes most commands in O(1) complexity and with minimal lines of code.

\*Redis on Windows:

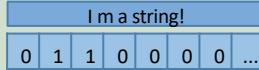
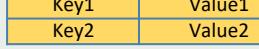
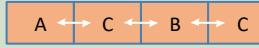
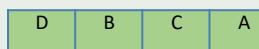
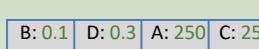
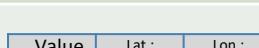
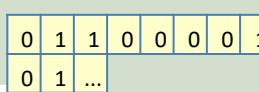
<https://github.com/MSOpenTech/redis>

# Redis Architecture



\*from <https://matt.sh/redis-architecture-diagram>

# Redis data types

Redis Data		Contains	Read/write ability
Type			
String	 0 1 1 0 0 0 0 ...	Binary-safe strings (up to 512 MB), Integers or Floating point values, Bitmaps.	Operate on the whole string, parts, increment/decrement the integers and floats, get/set bits by position.
Hash	 Key1 Value1 Key2 Value2	Unordered hash table of keys to string values	Add, fetch, or remove individual items by key, fetch the whole hash.
List	 A ← C ← B → C	Doubly linked list of strings	Push or pop items from both ends, trim based on offsets, read individual or multiple items, find or remove items by value.
Set	 D B C A	Unordered collection of unique strings	Add, fetch, or remove individual items, check membership, intersect, union, difference, fetch random items.
Sorted Set	 B: 0.1 D: 0.3 A: 250 C: 250	Ordered mapping of string members to floating-point scores, ordered by score	Add, fetch, or remove individual items, fetch items based on score ranges or member value.
Geospatial index	 Value Lat.: 20.63373 Lon.: -103.55328	Sorted set implementation using geospatial information as the score	Add, fetch or remove individual items, search by coordinates and radius, calculate distance.
HyperLogLog	 0 1 1 0 0 0 0 1 0 1 ...	Probabilistic data structure to count unique things using 12Kb of memory	Add individual or multiple items, get the cardinality.

\*Redis data types internals: <https://cs.brown.edu/courses/cs227/archives/2011/slides/mar07-redis.pdf>

# Redis Commands - Basic

Strings	
<b>Get/Set strings</b> <i>redis&gt; SET foo "hello!"</i> <i>OK</i> <i>redis&gt; GET foo</i> <i>"hello!"</i>	SET [key value] / GET [key] <b>O(1)</b>
<b>Increment numbers</b> <i>redis&gt; SET bar 223</i> <i>OK</i> <i>redis&gt; INCRBY bar 1000 (integer) 1223</i>	INCRBY [key increment] <b>O(1)</b>
<b>Get multiple keys at once</b> <i>redis&gt; MGET foo bar</i> <i>1. "hello!" 2. "1223"</i>	MGET [key key ...] <b>O(N) : N=# of keys.</b>
<b>Set multiple keys at once</b> <i>&gt; MSET foo "hello!" bar 1223</i> <i>OK</i>	MSET [key value key value ...] <b>O(N) : N=# of keys.</b>
<b>Get the length of a string</b> <i>redis&gt; STRLEN foo (integer) 6</i>	STRLEN [key] <b>O(1)</b>

<b>Update a value retrieving the old one</b> <i>redis&gt; GETSET foo "bye!"</i> <i>redis&gt; GET foo "bye!"</i>	GETSET [key value] <b>O(1)</b>
---	--------------------------------

# Redis Commands - Basic

Keys	
<b>Key removal</b>	DEL [key ...] <b>O(1)</b>
<i>redis&gt; DEL foo (integer) 1</i>	
<b>Test for existence</b>	EXISTS [key ...] <b>O(1)</b>
<i>redis&gt; EXISTS foo (integer) 1</i>	
<b>Get the type of a key</b>	TYPE [key] <b>O(1)</b>
<i>redis&gt; TYPE foo string</i>	
<b>Rename a key</b>	RENAME [key newkey] <b>O(1)</b>
<i>redis&gt; RENAME bar new_bar</i>	
<i>OK</i>	<b>O(1)</b>
<i>redis&gt; EXPIRE foo 10 (integer) 1</i>	
<b>Get key time-to-live</b>	TTL [key] <b>O(1)</b>
<i>redis&gt; TTL foo (integer)</i>	
<i>10</i>	

# Redis Commands - Lists & Hashes

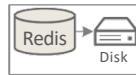
Lists		Hashes	
<b>Push on either end</b> <i>redis&gt; RPUSH jobs "foo"</i> <i>(integer) 1</i> <i>redis&gt; LPUSH jobs "bar"</i> <i>(integer) 1</i>	RPUSH/LPUSH [key value] <b>O(1)</b>	<b>Set a hashed value</b> <i>redis&gt; HSET user:1 name John (integer) 1</i>	HSET [key field value] <b>O(1)</b>
<b>Pop from either end</b> <i>redis&gt; RPOP jobs</i> <i>"foo"</i> <i>redis&gt; LPOP jobs</i> <i>"bar"</i>	RPOP/LPOP [key] <b>O(1)</b>	<b>Set multiple fields</b> <i>redis&gt; HMSET user:1 lastname Smith visits 1 OK</i>	HMSET [key field value ...] <b>O(1)</b>
<b>Blocking Pop</b> <i>redis&gt; BLPOP jobs</i> <i>redis&gt; BRPOP jobs</i>	BRPOP/BLPOP [key] <b>O(1)</b>	<b>Get a hashed value</b> <i>redis&gt; HGET user:1 name "John"</i>	HGET [key field] <b>O(1)</b>
<b>Pop and Push to another list</b>	RPOPLPUSH [src dst] <b>O(1)</b>	<b>Get all the values in a hash</b> <i>redis&gt; HGETALL user:1</i> 1) "name" 2) "John" 3) "lastname" 4) "Smith" 5) "visits" 6) "1"	HGETALL [key] <b>O(N)</b> : N=size of hash.
<b>Get an element by index</b> <i>redis&gt; LINDEX jobs 1</i> <i>"foo"</i>	LINDEX [key index] <b>O(N)</b>	<b>Increment a hashed value</b> <i>redis&gt; HINCRBY user:1 visits 1</i> <i>(integer) 2</i>	HINCRBY [key field incr] <b>O(1)</b>
<b>Get a range of elements</b> <i>redis&gt; LRANGE jobs 0 -1</i>	LRANGE [key start stop] <b>O(N)</b>		

# Redis Commands - Sets & Sorted sets

Sets	Sorted sets		
<b>Add member to a set</b> <code>redis&gt; SADD admins "Peter"</code> (integer) 1 <code>redis&gt; SADD users "John" "Peter"</code> (integer) 2	<code>SADD [key member ...]</code> $O(1)$	<b>Add member to a sorted set</b> <code>redis&gt; ZADD scores 100 "John"</code> (integer) 1 <code>redis&gt; ZADD scores 50 "Peter" 200 "Charles" 1000 "Mary"</code> (integer) 3	<code>ZADD [key score member]</code> $O(\log(N))$
<b>Pop a random element</b> <code>redis&gt; SPOP [key]</code> $O(1)$  <code>redis&gt; SPOP users "John"</code>		<b>Get the rank of a member</b> <code>redis&gt; ZRANK scores "Mary"</code> (integer) 3	<code>ZRANK [key member]</code> $O(\log(N))$
<b>Get all elements</b> <code>redis&gt; SMEMBERS users</code> 1) "Peter" 2) "John"	<code>SMEMBERS [key]</code> $O(N)$ : N=size of set.	<b>Get elements by score range</b> <code>redis&gt; ZRANGEBYSCORE scores 200 +inf WITHSCORES</code> 1) "Charles" 2) 200 3) "Mary" 4) 1000	<code>ZRANGEBYSCORE [key min max]</code> $O(\log(N))$
<b>Union multiple sets</b> <code>redis&gt; SUNION users admins</code> 1) "Peter" 2) "John"	<code>SUNION [key key ...]</code> $O(N)$	<b>Increment score of member</b> <code>redis&gt; ZINCRBY scores 10 "Mary" "1010"</code>	<code>ZINCRBY [key incr member]</code> $O(\log(N))$
<b>Diff. multiple sets</b> <code>redis&gt; SDIFF users admins</code> 1) "John"	<code>DIFF [key key ...]</code> $O(N)$	<b>Remove range by score</b> <code>redis&gt; ZREMRANGEBYSCORE scores 0 100</code> (integer) 2	<code>ZREMRANGEBYSCORE [key min max]</code> $O(\log(N))$

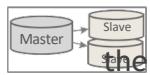
# Scaling Redis

## ► Persistence



Redis provides two mechanisms to deal with persistence: Redis database snapshots (RDB) and append-only files (AOF).

## ► Replication



A Redis instance known as the *master*, ensures that one or more instances known as *slaves*, exact copies of the master. Clients can connect to the master or to the slaves. Slaves are read only by default.

## ► Partitioning

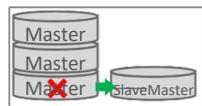


Breaking up data and distributing it across different hosts in a cluster.

Can be implemented in different layers:

- ▶ **Client:** Partitioning on client-side code.
- ▶ **Proxy:** An extra layer that proxies all redis queries and performs partitioning (i.e. [Twemproxy](#)).
- ▶ **Query Router:** instances will make sure to forward the query to the right node. (i.e [Redis Cluster](#)).

## ► Failover



- ▶ Manual
- ▶ Automatic with Redis Sentinel (for master-slave topology)
- ▶ Automatic with Redis Cluster (for cluster topology)

# Redis topologies

I.

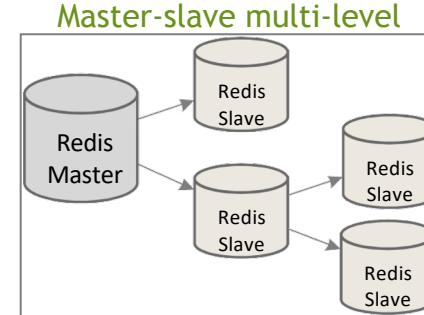
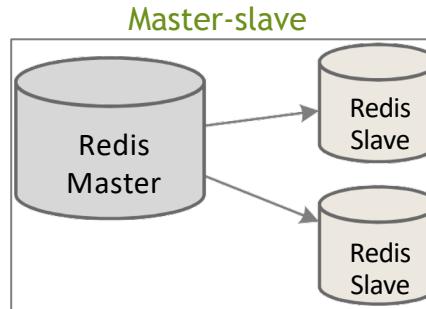
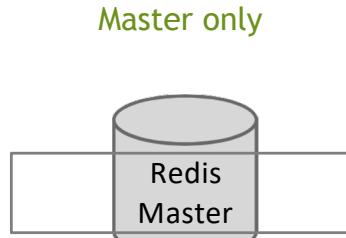
II. Sentinel      Standard  
alone

III. Twemproxy      (distribute data)  
y                      (automatic failover and distribute data)

IV. Cluster

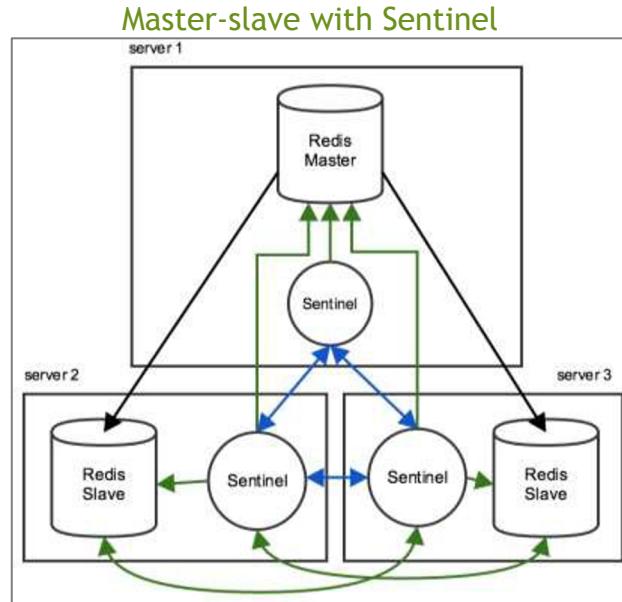
# Redis topologies I - Standalone

- The master data is optionally replicated to slaves.
- The slaves provides data redundancy, reads offloading and save-to-disk offloading.
- Clients can connect to the Master for read/write operations or to the Slaves for read operations.
- Slaves can also replicate to its own slaves.
- There is no automatic failover.



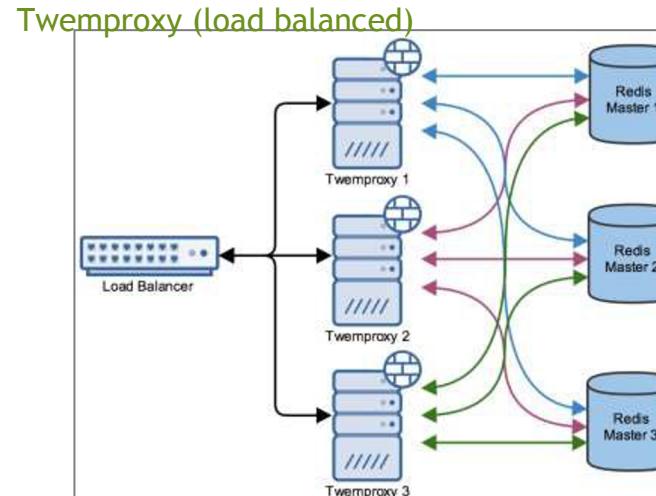
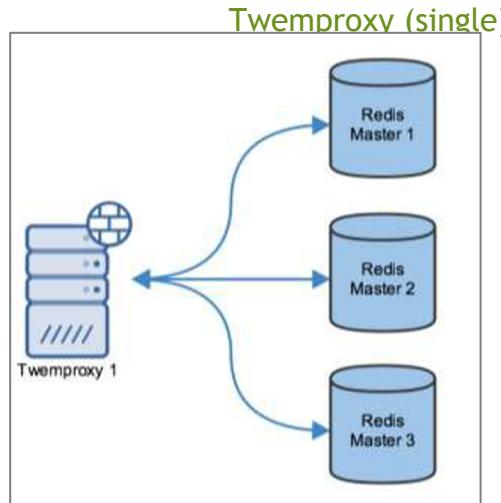
# Redis topologies II - Sentinel

- ▶ Redis Sentinel provides a reliable **automatic failover** in a master/slave topology, automatically promoting a slave to master if the existing master fails.
- ▶ Sentinel does not distribute data across nodes.



# Redis topologies III - Twemproxy

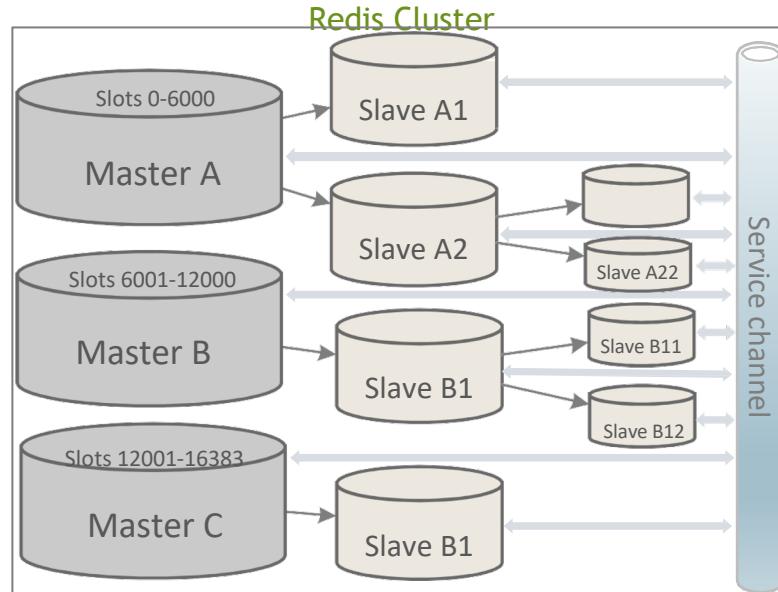
- ▶ Twemproxy\* works as a proxy between the clients and many Redis instances.
- ▶ Is able to automatically distribute data among different standalone Redis instances.
- ▶ Supports consistent hashing with different strategies and hashing functions.
- ▶ Multi-key commands and transactions are not supported.



\*Twemproxy is a project from Twitter and is not part of redis:  
<https://github.com/twitter/twemproxy>

# Redis topologies IV - Cluster

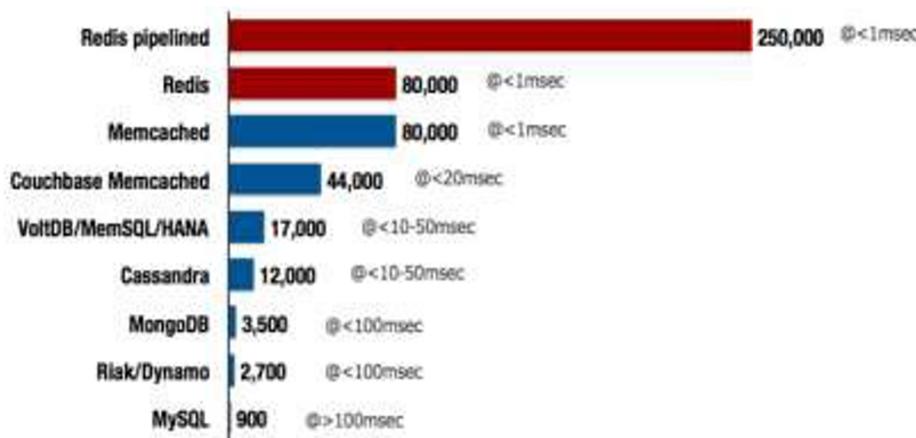
- ▶ Redis Cluster distributes data across different Redis instances and perform automatic failover if any problem happens to any master instance.
- ▶ All nodes are directly connected with a service channel.
- ▶ The keyspace is divided into hash slots. Different nodes will hold a subset of hash slots.
- ▶ Multi-key commands are only allowed for keys in the same hash slot.



# Redis advantages

- ▶ Performance
- ▶ Availability
- ▶ Fault-Tolerance
- ▶ Scalability (adaptability)
- ▶ Portability

- ▶ Support for complex data types and structures
- ▶ Atomic operations and Transactions
- ▶ Pipelining (send multiple commands at once)
- ▶ LUA scripting support
- ▶ LRU eviction of keys
- ▶ Keys with limited time-to-live
- ▶ Simple to install, setup and manage
- ▶ Highly configurable
- ▶ Supported on many languages
- ▶ Straightforward and well documented API
- ▶ Open Source
- ▶ Available on Azure



NoSQL & SQL response performance comparison  
(from <https://redislabs.com/blog/the-proven-redis-performance>)

# Key-value and Caching

# What is caching?

- From Wikipedia: "A cache is a collection of data duplicating original values stored elsewhere or computed earlier, where the original data is expensive to fetch (owing to longer access time) or to compute, compared to the cost of reading the cache."
- Term introduced by IBM in the 60's

# The anatomy

- simple key/value storage
- simple operations
  - save
  - get
  - delete

# Terminology

- **storage cost**
- **retrieval cost** (network load / algorithm load)
- **invalidation** (keeping data up to date / removing irrelevant data)
- **replacement policy** (FIFO/LFU/LRU/MRU/RANDOM vs. Belady's algorithm)
- **cold cache / warm cache**

# Terminology

- cache hit and cache miss
- typical stats:
  - hit ratio (hits / hits + misses)
  - miss ratio (1 - hit ratio)
- 45 cache hits and 10 cache misses
  - $45/(45+10) = 82\%$  hit ratio
  - 18% miss ratio

# When to cache?

- caches are only efficient when the benefits of faster access outweigh the overhead of checking and keeping your cache up to date
  - more cache hits than cache misses

# Where are caches used?

- at hardware level (cpu, hdd)
- operating systems (ram)
- web stack
- applications
- your own short term vs long term memory

# Caches in the web stack

- Browser cache
- DNS cache
- Content Delivery Networks (CDN)
- Proxy servers
- Application level
  - **full output caching** (eg. Wordpress WP-Cache plugin)
- ...

# Efficiency of caching?

- the earlier in the process, the closer to the original request(er), the faster
  - browser cache will be faster than cache on a proxy
- but probably also the harder to get it right
  - the closer to the requester the more parameters the cache depends on



# Memcached

# About memcached

- Free & open source, high-performance, distributed memory object caching system
- Generic in nature, intended for use in speeding up dynamic web applications by alleviating database load.
- key/value dictionary

# About memcached

- Developed by Brad Fitzpatrick for LiveJournal in 2003
- Now used by Netlog, Facebook, Flickr, Wikipedia, Twitter, YouTube ...

# Technically

- It's a server
- Client access over TCP or UDP
- Servers can run in pools
  - eg. 3 servers with 64GB mem each give you a single pool of 192GB storage for caching
- Servers are independent, clients manage the pool

# What to store in memcache?

- **high demand** (used often)
- **expensive** (hard to compute)
- **common** (shared accross users)
- Best: All three

# What to store in memcache? (cont'd)

- Typical:
  - user sessions (often)
  - user data (often, shared)
  - homepage data (eg. often, shared, expensive)

# Memcached principles

- **Fast network access** (memcached servers close to other application servers)
- **No persistency** (if your server goes down, data in memcached is gone)
- **No redundancy / fail-over**
- **No replication** (single item in cache lives on one server only)
- **No authentication** (not in shared environments)

# Memcached principles (cont'd)

- **1 key is maximum 1MB**
- **keys are strings of 250 characters** (in application typically MD5 of user readable string)
- **No enumeration of keys** (thus no list of valid keys in cache at certain moment)
- **No active clean-up** (only clean up when more space needed, LRU: Least Recently Used )

# PHP Client functions

**Memcached::add** — Add an item under a new key

**Memcached::addServer** — Add a server to the server pool

**Memcached::decrement** — Decrement numeric item's value

**Memcached::delete** — Delete an item

**Memcached::flush** — Invalidate all items in the cache

**Memcached::get** — Retrieve an item

**Memcached::getMulti** — Retrieve multiple items

**Memcached::getStats** — Get server pool statistics

**Memcached::increment** — Increment numeric item's value

**Memcached::set** — Store an item

```
<?php

$html = $cache->get('mypage');
if (!$html)
{
    ob_start();
    echo "<html>";
    //      all the      stuff      goes here
    echo "</html>";
    $html = ob_get_contents();
    ob_end_clean();
    $cache->set('mypage',      $html);
}
echo $html;

?>
```

# Data caching

- on a lower level
- easier to find all dependencies
- ideal solution for offloading database queries
- the database is almost always the biggest bottleneck in backend performance problems

```
<?php

function    getUserData($UID)
{
    $key =  'user_' . $UID;
    $userData   = $cache->get($key);
    if (!$userData)
    {
        $queryResult    = Database::query("SELECT      * FROM USER
WHERE uid = " . (int) $UID);
        $userData   = $queryResult->getRow();
        $cache->set($userData);
    }
    return $userData;
}

?>
```

“There are only two  
hard things in  
Computer Science:  
cache invalidation  
and naming things.”

Phil Karlton

# Invalidation

- Caching for a certain amount of time
  - eg. 10 minutes
  - don't delete caches
  - thus: You can't trust that data coming from cache is correct

# Invalidation (cont'd)

- Use: Great for summaries
  - Overview
  - Pages where it's not that big a problem if data is a little bit out of date (eg. search results)
- Good for quick and dirty optimizations

# Invalidation (cont'd)

- Store forever, and expire on certain events
  - the userdata example
  - store userdata for ever
  - when user changes any of his preferences, throw cache away

```
<?php
function      getUserData($UID,      $invalidateCache      =  false)
{
    $db = DB::getInstance();
    $db->prepare("SELECT *
                  FROM
                  USERS      {UID}");
    $db->assignInt(1,$uid =  $UID);
    $db->setCacheTTL(0); // if      cache forever
    ($invalidateCache)
    {
        return  $db->invalidateCache();
    }
    $db->execute(); return $db-
    >getRow();
}
?>
```

```
<?php
function updateUserData($UID,      $data)
{
    $db = DB::getInstance();
    $db->prepare("UPDATE USERS
                  SET ...
                  WHERE uid = {UID}");
```

...

```
getUserData($UID,      true); // cache
                                invalid
return $result;               ate
}
?>
```

# Multi-Get Optimisations

- We reduced database access
- Memcached is faster, but access to memcache still has it's price
- Solution: multiget
  - fetch multiple keys from memcached in one single call
  - result is array of items

# More tips ...

- **Be carefull when security matters.**  
(Remember ‘no authentication’?)
  - Working on authentication for memcached via SASL Auth Protocol
- **Caching is not an excuse not to do database tuning.** (Remember cold cache?)
- **Make sure to write unit tests for your caching classes and places where you use it.**  
(Debugging problems related to out-of-date cache data is hard and boring. Very boring.)



POLITECNICO  
MILANO 1863

SYSTEMS AND METHODS FOR BIG AND UNSTRUCTURED DATA

# Key-value Databases – Redis

Marco Brambilla

[marco.brambilla@polimi.it](mailto:marco.brambilla@polimi.it)

 @marcobrambi