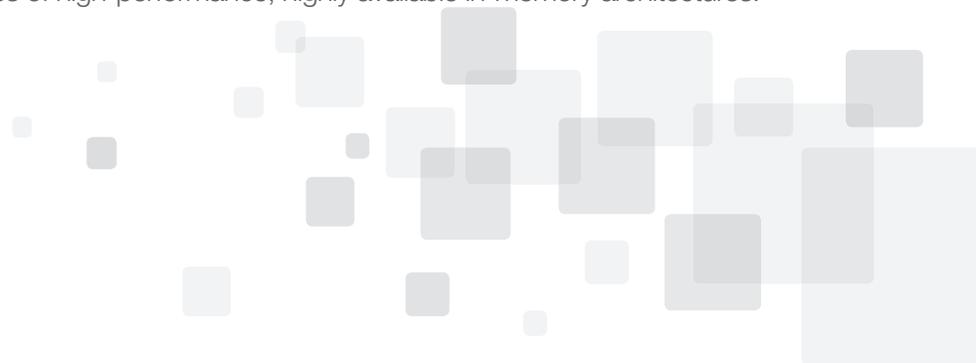


# Ditch the Disk:

## Designing a High-Performance In-Memory Architecture

The need for taking real-time action on Big Data intelligence is driving big changes to the traditional enterprise architecture. In particular, enterprises are increasingly keeping data instantly available in ultra-fast machine memory, rather than locking it away in slow, disk-bound databases.

In this paper, we explain why an architecture built around in-memory data management is the best model for achieving high performance and extremely low, predictable latency at scale. We examine why traditional, disk-based database systems are simply no longer sufficient for emerging real-time Big Data applications, and we enumerate the most important attributes of high-performance, highly available in-memory architectures.



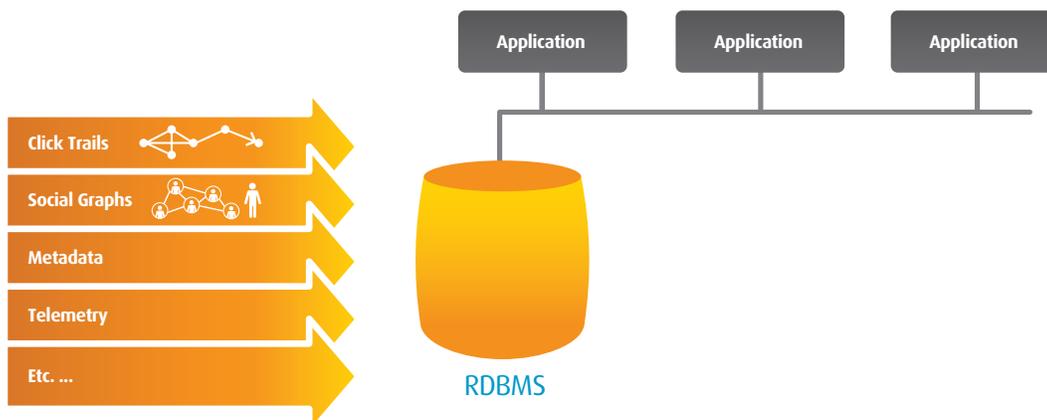
## Why It's Time to Ditch the Disk

Until recently, the architecture of enterprise computing had not changed for decades. Visit any corporate data center today, and you'll see pretty much the same thing: applications depositing their data into disk-bound relational databases. When those applications need that data again, they ask the disks for a copy, do something with it, and update the disk when necessary.

The standard architecture worked until the dawn of the Internet age. Back then, most business data fit easily in a single database server and could credibly be coerced into a relational format. In addition, the neat separation of concerns between application logic and data management gave rise to a convenient differentiation between application programmers on the one hand, and database programmers and administrators on the other.

In the early 2000s, however, cracks began to appear in the standard architecture. Exposed to growing Internet usage and sudden traffic spikes, applications and databases were more and more frequently brought to their knees. Enterprises struggled to build scalability into nearly every aspect of their infrastructure. First steps usually included scaling up by adding bigger, faster hardware. Then, when the practical limits of hardware scaling were reached, it was common practice to try to take advantage of parallelism to distribute load across many hardware instances.

However, for applications that relied on back-end relational database management systems (RDBMS) as their disk-based data hubs, the database always became a bottleneck to performance and scale. The speed and data volume requirements of Big Data only amplify these challenges, causing more and more enterprises to look for an alternative to the traditional disk-based paradigm.



**The new data you need to take action on comes in all shapes and sizes, from many sources. and no longer fits**

Figure 1: Trying to fit 21<sup>st</sup> century data into 20<sup>th</sup> century architecture

Businesses increasingly want to manipulate data, the shape and scale of which no longer fits neatly into the old architectural paradigm. Under the old paradigm (Figure 1), new data streams must be converted to a structured relational format and imported into the backend RDBMS. However, the majority of business information—by some estimates as much as 80%<sup>1</sup>—is unstructured data, which traditional, structured databases are ill-equipped to handle. As many as 10 to 20 percent of RDBMS implementations are attempts to shoe-horn unstructured data into the old structured-data-only model.<sup>2</sup> The expense and complexity of coercing new data streams into the old model is simply unsustainable.

In the absence of a suitable modern data architecture, application builders are forced to create ad-hoc data handling facilities for each application (Figure 2). While often a point-in-time necessity, such attempts lead to inconsistent speed, scale, and reliability. Stuffing data management concerns into business applications makes those applications more complex to build, riskier to manage, and more expensive to maintain.

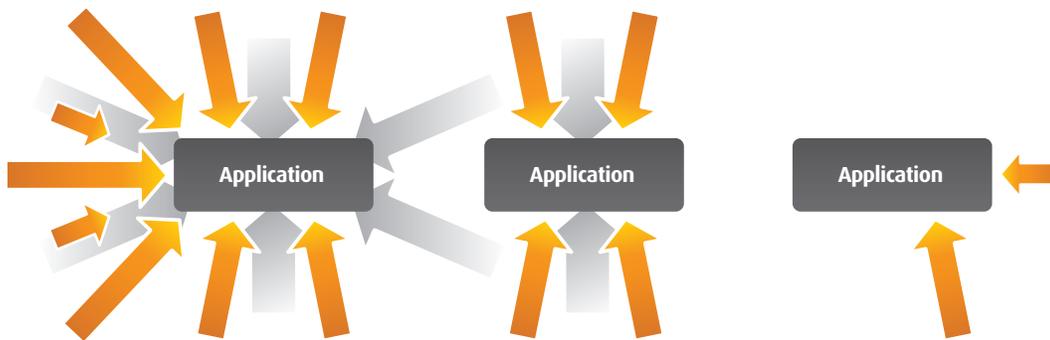


Figure 2: Building ad-hoc data management into each application leads to inconsistent speed, scale, and reliability

Today's enterprises are coming to understand that the traditional disk-bound, database-centric architecture needs to be replaced by something new. It will no longer be possible to force enterprise customers, employees, and partners to wait while applications deal with disks. Instead, to achieve high performance and extremely low, predictable latency at scale, enterprises are moving data into ultra-fast machine memory.

<sup>1</sup> Bank of America Merrill Lynch, Big Data II, March 2012, p. 21

<sup>2</sup> *Ibid*

## In-Memory Data Management:

### ACHIEVING PERFORMANCE AT BIG DATA SCALE

To take full advantage of Big Data, enterprises must adopt a new data management architecture. It has to be fast and scalable, ready to handle new data in all shapes and sizes. It must also possess the virtues of long-established data management systems, including consistency, durability, high availability, monitoring and management. Finally, those capabilities must be available to all applications across the enterprise, its partner networks, and its customer ecosystems.

The answer, of course, is shifting data from traditional databases to modern, ultra-fast in-memory data stores. In-memory data management slashes application response times from milliseconds to microseconds because data is made available where computation happens, optimizing existing operations and enabling entirely new kinds applications (Figure 3).

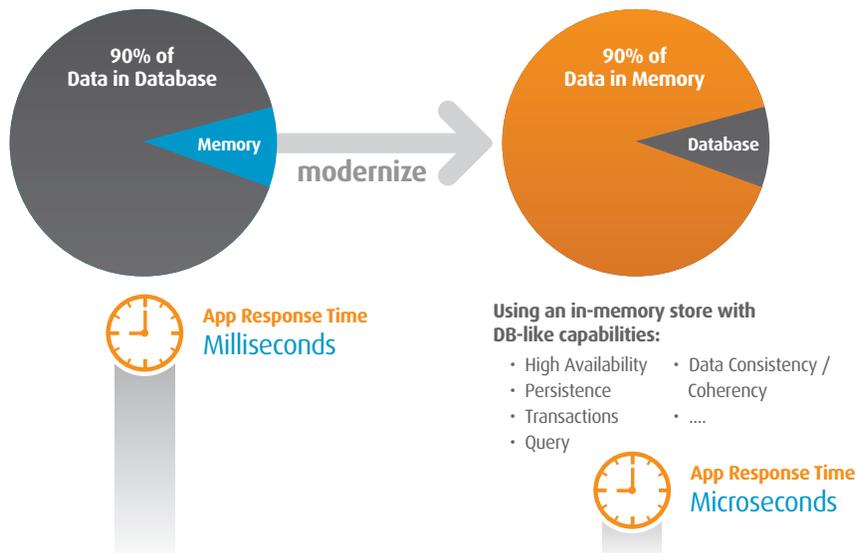


Figure 3: Moving data in memory transforms what's possible

In-memory computing is powering innovative disruption in a range of industries. Here are just a few examples:

- **Financial services:** Risk management and fraud detection happen in real time, during a transaction, rather than hours later.
- **E-commerce:** With entire product catalogs, shopping carts, and transaction histories in memory, retailers can deliver faster online commerce experiences, including real-time cross-sell offers.
- **Media and entertainment:** By moving Big Data-sized sets of profile data, content metadata, and rules into memory, media and online gaming companies can deliver richer streaming experiences and gameplay.
- **Logistics:** Shipping, transport, and tracking companies are managing geolocation and inventory data in RAM, delivering real-time route optimization.
- **Hospitality and Travel:** Hotels, airlines, and online travel sites are putting entire reservation systems into memory, speeding customer interactions and allowing for real-time yield management.
- **Government:** Law enforcement and other agencies are moving crucial data into memory to gather and share it in real time.
- **Marketing Services:** Advertising-tracking and market data providers are using in-memory computing to deliver instant access to ad and product performance data.

## The 6 Most Crucial Attributes of a High-Performance In-Memory Architecture

Of course, while local memory is very fast, it is also volatile. If not architected properly, a scaled-out application's in-memory data can easily become inconsistent across instances of the application. The move from disk-based to memory-based data architectures requires a robust in-memory data management architecture that delivers high speed, low-latency access to terabytes of data across a distributed application, while maintaining capabilities previously provided by the RDBMS such as data consistency, durability, high availability, fault tolerance, monitoring, and management.

Here are six of the most important concerns to address when evaluating in-memory data management solutions:

### 1. Predictable, Extremely Low Latency

Working with data in machine memory is orders of magnitude faster than moving it over a network or getting it from a disk. This speed advantage is critical for real-time data processing at the scale of Big Data. However, Java garbage collection is an Achilles' heel when it comes to using large amounts of in-memory data. While terabytes of RAM are available on today's commodity servers, Java applications can only use a few gigabytes of that before long, unpredictable garbage collection pauses cause application slowdowns and violations of service-level agreements (SLAs).

For Java applications to exhibit low and predictable latency, they need an in-memory management solution that can manage terabytes of data without suffering from garbage collection pauses.

### 2. Easy Scaling with Minimal Server Footprint

Scaling to terabytes of in-memory data should be easy—and shouldn't require the cost and complexity of dozens of servers and hundreds of JVMs. Your in-memory management solution should be able to scale up as much as possible on each machine so that you're not saddled with managing and monitoring a 100-node data grid. By fully utilizing the RAM on each server, you can dramatically reduce not only hardware costs, but also personnel costs associated with monitoring large server networks.

### 3. Fault Tolerance and High Availability

Mission-critical applications demand fault tolerance and high availability. The volatile nature of in-memory data requires a data management solution that delivers five-nines (99.999%) uptime with no data loss and no single points of failure.

### 4. Distributed In-Memory Stores with Data Consistency Guarantees

With the rise of in-memory data management as a crucial piece of Big Data architectures, businesses increasingly rely on having tens of terabytes of data accessible for real-time, mission-critical decisions. Multiple applications (and instances of those applications) will need to tap in-memory stores that are distributed across multiple servers. Thus, in-memory architectures must ensure the consistency and durability of critical data across that array. Ideally, you'll have flexibility in choosing the appropriate level of consistency guarantees, from eventual and strong consistency up to transactional consistency.

### 5. Fast Restartability

In-memory architectures must allow for quickly bringing machines back online after maintenance or other outages. Systems designed to backup and restore only a few gigabytes of in-memory data often exhibit pathological behavior around startup, backup, and restore as data sizes grow much larger. In particular, recreating a terabyte-sized in-memory store can take days if fast restartability is not a tested feature. Hundreds of terabytes? Make that weeks.

## 6. Advanced In-Memory Monitoring and Management Tools

In dynamic, large-scale application deployments, visibility and management capabilities are critical to optimizing performance and reacting to changing conditions. Control over where critical data is and how it is accessed by application instances gives operators the edge they need to anticipate and respond to significant events like load spikes, I/O bottlenecks or network and hardware failures before they become problems. Your in-memory architecture should be supplemented with a clear dashboard for understanding up-to-the-millisecond performance of in-memory stores, along with easy-to-use tools for configuring in-memory data sets.



## New Ad Platform Taps BigMemory for Unprecedented Speed at Massive Scale

---

**“I WANTED TO THROW OUT THE DATABASE AND, WITH IT, THE DISKS. WITH BIGMEMORY, I CAN.”**

– BEN LINDQUIST, VP OF TECHNOLOGY, ADJUGGLER

---

### THE BIG CHALLENGE: HANDLE MASSIVE TRANSACTION VOLUMES WITH SPEED AND RELIABILITY

AdJuggler is a global leader in on-demand digital ad management technology and media services. The company enables publishers and networks to source and manage all media through an integrated platform that maximizes yield based on the value and preferences of each unique user. AdJuggler's online advertising server is essentially an accounting system -- but with a huge transaction volume that generates terabytes of data every day, making reliability, timeliness of data, and speed at scale tremendously important. To support clients with fluid buying and selling of ad inventory, AdJuggler needed to scale far beyond the point of traditional business object storage.

### WHY BIGMEMORY: PROVIDES SIMPLE, ONE-VENDOR SOLUTION FOR COMPLEX ARCHITECTURAL DEMANDS

Ben Lindquist, VP of Technology at AdJuggler, says disk-bound database architectures are simply not up to the demands of his task. In place of a disk-bound architecture, Lindquist deployed Terracotta BigMemory Max as AdJuggler's distributed, in-memory system of record. “I wanted to throw out the database and, with it, the disks,” Lindquist says. “At AdJuggler, we're building a one million transaction-per-second online advertising marketplace. Speed at scale is everything, and we are past the point where we can do things in traditional ways. Terracotta and BigMemory will allow us to scale the backend of our system beyond a million transactions per second and make it globally distributed. Disks are now relegated to archiving. With In-Genius we will be able to utilize its complex event processing to automatically generate our accounting records, giving us a distinct advantage over our competition.”

### THE BIG RESULTS: SHARED STATE TRANSACTIONS AT LOW LATENCY OPTIMIZES ONLINE AD CAMPAIGNS; ENSURES INCOME FOR ALL TRANSACTIONS

For AdJuggler, BigMemory Max has become a critical infrastructure component that is the company's system of record for shared state. Whether sharing state between east and west coasts of the United States or with Europe or Asia, BigMemory Max allows AdJuggler to optimize online ad campaigns by enabling its transaction engine to make updates and replicate them. “The ability to have a low-latency update of the shared state has a direct impact on the bottom line,” says Lindquist.

Next, Lindquist wants to use Terracotta's BigMemory-Hadoop Connector to make Hadoop insights available to the transaction engine in real time, arming the Adjugger platform with up-to-the-millisecond intelligence. "We plan to use BigMemory as a hub for Hadoop," Lindquist explains. "We'll pull data into Hadoop directly out of BigMemory, and then feed Hadoop results right back in."

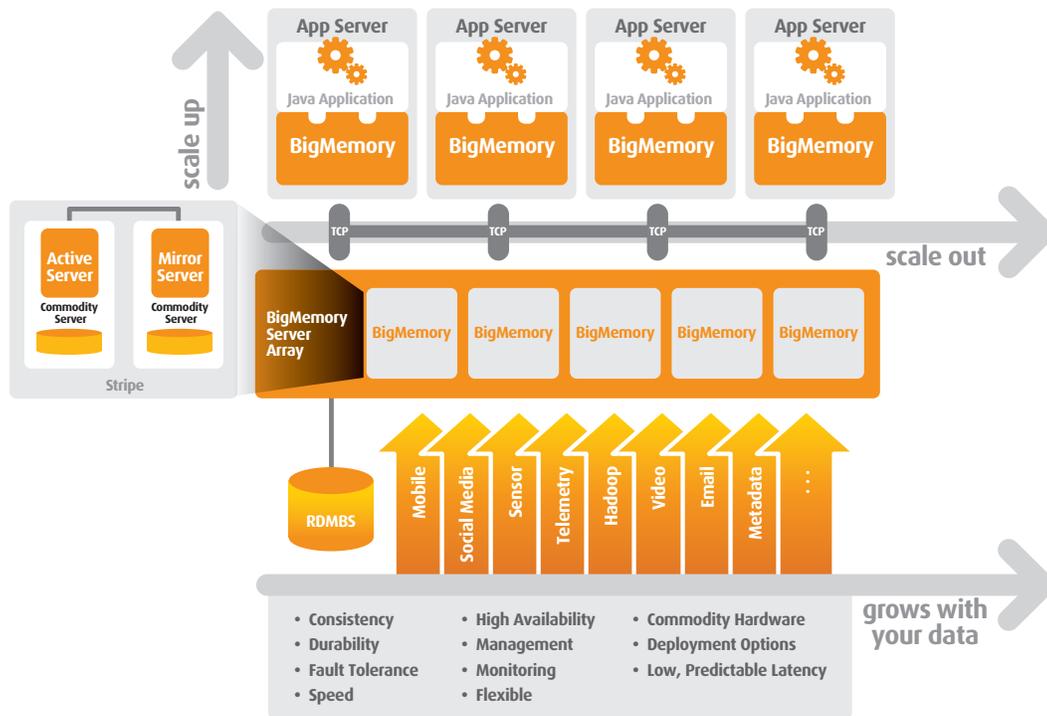


Figure 4: Ad Juggler is using BigMemory Max as a hub for its next-generation ad server

## Ready to Ditch Your Disks?

Terracotta software is already used in millions of deployments worldwide to make petabytes of data available in machine memory at microsecond speed. As big data goes mainstream, Hadoop is finding a home in many of those same IT shops. BigMemory's snap-in integration with Hadoop makes real-time intelligence powered by big data a reality for the enterprise.

### FOR MORE INFORMATION:

Email [sales@terracotta.org](mailto:sales@terracotta.org) or download BigMemory for free at <http://terracotta.org/bigmemory>

## About Terracotta, Inc.

Terracotta, the leader in in-memory technologies for enterprise Big Data, is the innovator behind some of the most widely used software for application scalability, availability and performance. Headquartered in San Francisco, Terracotta serves the majority of Global 2000 companies as customers and boasts more than 2.5 million software installations worldwide. The company's flagship BigMemory platform is an in-memory data management solution delivering performance at Big Data scale. Terracotta's other leading solutions include Ehcache, the award-winning de facto caching standard for enterprise Java, and Quartz Scheduler, a leading job scheduler. Terracotta is a wholly-owned subsidiary of Software AG (Frankfurt TecDAX: SOW). For more information, visit [www.terracotta.org](http://www.terracotta.org) or follow Terracotta on Twitter, Facebook and LinkedIn.

**For more information, please visit [www.terracotta.org](http://www.terracotta.org).**

### **TERRACOTTA, INC.**

575 Florida St. Suite 100  
San Francisco, CA 94110

For Product, Support, Training and Sales Information:  
[sales@terracotta.org](mailto:sales@terracotta.org)

#### **USA Toll Free**

+1-888-30-TERRA

#### **International**

+1-415-738-4000

#### **Terracotta China**

[china@terracottatech.com](mailto:china@terracottatech.com)  
+1-415-738-4088

