#### **OPEN SOURCE BUSINESS CONFERENCE**

E ald as

**Building Your Big Data Future with Open Source** 

# COMPUTERWORLD OSBC SAN FRANCISCO



Matthew Aslett, The 451 Group matthew.aslett@the451group.com

© 2011 by The 451 Group. All rights reserved



M Cho and

#### **Overview**

- Open source impact in the database market
- NoSQL and NewSQL databases
  - Adoption/development drivers
  - Use cases
- And beyond
  - Data grid/cache
  - Cloud databases
  - Total data



**COMPUTERWORLD** 

451 Research is focused on the business of enterprise IT innovation. The company's analysts provide critical and timely insight into the competitive dynamics of innovation in emerging technology segments.



Tier1 Research is a single-source research and advisory firm covering the multi-tenant datacenter, hosting, IT and cloud-computing sectors, blending the best of industry and financial research.

#### **Uptime**Institute

The Uptime Institute is 'The Global Data Center Authority' and a pioneer in the creation and facilitation of end-user knowledge communities to improve reliability and uninterruptible availability in datacenter facilities.



TheInfoPro is a leading IT advisory and research firm that provides real-world perspectives on the customer and market dynamics of the enterprise information technology landscape, harnessing the collective knowledge and insight of leading IT organizations worldwide.



ChangeWave Research is a research firm that identifies and quantifies 'change' in consumer spending behavior, corporate purchasing, and industry, company and technology trends.



de a

#### **Coverage areas**

- Senior analyst, enterprise software
- Commercial Adoption of Open Source (CAOS)
  - Adoption by enterprises
  - Adoption by vendors
- Information Management
  - Databases
  - Data warehousing
  - Data caching

#### **Relevant reports**

Turning the Tables

COMPUTERWORLD

- The impact of open source on the enterprise database market
- Published March 2008
- Survey of 368 database purchasers
- sales@the451group.com

Turning the Tables? The impact of open source on the enterprise database market

451 COMMERCIAL ADOPTION OF OPEN SOURCE (CAOS) RESEARCH SERVICE

REPORT 7, MARCH 2008







#### **Relevant reports**

- Warehouse Optimization

COMPUTERWORLD

SAN FRANCISCO

- Ten considerations for choosing/building a data warehouse
- Published September 2009
- The role of open source and the emergence of Hadoop
- sales@the451group.com



The second second

(Same Alexand



#### **Relevant reports**

- Data Warehousing 2009-2013
  - Market Sizing, Landscape and Future
  - Published August 2010

COMPUTERWORLD

- The potential impact of Hadoop
- sales@the451group.com



20.

#### DATA WAREHOUSING: 2009-2013 Meter Sizing, Landrages and Fature

This report provides market-sizing outimates for the data-warehouse sector from 2000 to 2013. It includes revenue estimates and growth projections, and examines the husiness and technology irends driving this market.

#### IM INFORMATION MANAGEMENT

4 FINDINGS	5 IMPLICATIONS	4 BOTTOM LINE
<ul> <li>FOLT Withins dominate the data superioran market, with USUE.</li> <li>of latal researce in 2010. These sensities are expected in mide their schedulings and generate INVEST of measure in 2013. PAGE 6</li> </ul>	Protock that late advantage of ingeneral backware performance will obtain researce growth for all results, and will protock the market share of insurances, PROE 40     Au a result of spatience performance.	<ul> <li>We estimate that her data workstaating market will not a compared around growth rate of 11.12% here 2009 Hereigh 2013 in market is init of \$11.2% in revenues.</li> </ul>
<ul> <li>Analytic databases are new slife to take advantage of greater processor performance at a lower real, improving price/performance and</li> </ul>	ingenerating, date standarding southers are date bidling obtaining of the apparticity in bring more advanced analytic separabilities in the Dil wogles. PAGE 44	
<ul> <li>Imming lawrine in may, PABE 60</li> <li>The invest-driving adaptive of data monochronology frameri i dravgoli stave tier mark dapa of the indexedagy papelasite includences in many angledatie technicases in many angle PAGE 8</li> <li>We be application of shard applications of an of the index papelasites, and some for paratient of paths of monophies fully</li> </ul>	<ul> <li>All mapping memory and the second seco</li></ul>	
ikai many antikal managemeni wito distrikatel nar and materi. PACE 45	<ul> <li>Webs the Hashamy Gare is not a silent alternative in traditional analysis CDs, the increased materia of associated projects means that our sames for Hashamy and MagDestare resulted analysis CDs will see lays. PAGE 49</li> </ul>	ALIGUST 2010
	THE GO	GELP INCOMING HARASHINT





#### **Relevant reports**

#### - NoSQL, NewSQL and Beyond

- Assessing the drivers behind the development and adoption of NoSQL and NewSQL databases, as well as data grid/caching
- Released April 2011

COMPUTERWORLD

SAN FRANCISCO

- Role of open source in driving innovation
- sales@the451group.com



# NoSQL, NewSQL and Beyond

#### – NoSQL

COMPLETERWORL

- New non-relational databases
- Rejection of fixed table schema and join operations
- Meet scalability requirements of distributed architectures
- And/or schema-less data management requirements
- Data grid/cache

ശNewSQL

- New relational databases
- Retain SQL and ACID compliance
- Meet scalability requirements of distributed architectures
- Or to improve performance such that horizontal scalability is no longer a necessity
- Store data in memory to increase app and database performance
- Potential alternative to relational databases as the primary platform for distributed data management

# NoSQL, NewSQL and Beyond

#### – NoSQL

COMPLETERWORL

- Big tables
- Key value stores
- Document stores
- Graph databases

**C3NewSQL** 

- MySQL storage engines
- Transparent sharding
- Appliances: software and hardware
- New databases

#### Data grid/cache

- Spectrum of data management capabilities
- From non-persistent data caching to persistent caching, replication, and distributed data and compute grid functionality







#### **Database SPRAIN**

 "An injury to ligaments... caused by being stretched beyond normal capacity"

Wikipedia

M Cho in

- Six key drivers for NoSQL/NewSQL/DDG adoption
  - Scalability
  - Performance
  - Relaxed consistency
  - Agility
  - Intricacy
  - Necessity

# Scalability

- Associated sub-driver: Hardware economics
  - Scale-out across clusters of commodity servers
- Example project/service/vendor
  - BigTable, HBase, Riak, MongoDB, Couchbase, Hadoop
  - Amazon RDS, Xeround, SQL Azure, NimbusDB
  - Data grid/cache

COMPLETERWORL

- Associated use case:
  - Large-scale distributed data storage
  - Analysis of continuously updated data
  - Multi-tenant PaaS data layer



M the set

# Scalability

- User: StumbleUpon
- Problem:
  - Scaling problems with recommendation engine on MySQL
- Solution: HBase
  - Started using Apache HBase to provide real-time analytics on Su.pr
  - MySQL lacked the performance headroom and scale
  - Multiple benefits including avoiding declaring schema
  - Enables the data to be used for multiple applications and use cases

#### Performance

- Associated sub-driver: MySQL limitations
  - Inability to perform consistently at scale
- Example project/service/vendor
  - Hypertable, Couchbase, Membrain, MongoDB, Redis
  - Data grid/cache

COMPLETERWORL

- VoltDB, Clustrix
- Associated use case:
  - Real time data processing of mixed read/write workloads
  - Data caching
  - Large-scale data ingestion



#### Performance

- User: AOL Advertising
- Problem:
  - Real-time data processing to support targeted advertising
- Solution: Membase Server
  - Segmentation analysis runs in CDH, results passed into Membase
  - Make use of its sub-millisecond data delivery
  - More time for analysis as part of a 40ms targeted ad response time
  - Also real time log and event management



#### **Relaxed consistency**

- Associated sub-driver: CAP theorem
  - The need to relax consistency in order to maintain availability
- Example project/service/vendor
  - Dynamo, Voldemort, Cassandra
  - Amazon SimpleDB
- Associated use case:
  - Multi-data center replication
  - Service availability
  - Non-transactional data off-load



### **Relaxed consistency**

- User: Wordnik
- Problem:
  - MySQL too consistent –blocked access to data during inserts and created numerous temp files to stay consistent.
- Solution: MongoDB
  - Single word definition contains multiple data items from various sources
  - MongoDB stores data as a complete document
  - Reduced the complexity of data storage

# Agility

- Associated sub-driver: Polyglot persistence
  - Choose most appropriate storage technology for app in development
- Example project/service/vendor
  - MongoDB, CouchDB, Cassandra
  - Google App Engine, SimpleDB, SQL Azure
- Associated use case:
  - Mobile/remote device synchronization
  - Agile development
  - Data caching

COMPUTERWORL



M Cho all

# Agility

- User: Dimagi BHOMA (Better Health Outcomes through Mentoring and Assessments) project
- Problem:
  - Deliver patient information to clinics despite a lack of reliable Internet connections
- Solution: Apache CouchDB
  - Replicates data from regional to national database
  - When Internet connection, and power, is available
  - Upload patient data from cell phones to local clinic

# Intricacy

- Associated sub-driver: Big data, total data
  - Rising data volume, variety and velocity
- Example project/service/vendor
  - Neo4j, GraphDB, InfiniteGraph
  - Apache Cassandra, Hadoop,
  - VoltDB, Clustrix

COMPUTERWORL

- Associated use case:
  - Social networking applications
  - Geo-locational applications
  - Configuration management database



# Intricacy

- User: Evident Software
- Problem:
  - Mapping infrastructure dependencies for application performance management
- Solution: Neo4j
  - Apache Cassandra stores performance data
  - Neo4j used to map the correlations between different elements
  - Enables users to follow relationships between resources while investigating issues

# Necessity

- Associated sub-driver: Open source
  - The failure of existing suppliers to address the performance, scalability and flexibility requirements of large-scale data processing
- Example project/service/vendor
  - BigTable, Dynamo, MapReduce, Memcached
  - Hadoop, HBase, Hypertable, Cassandra, Membase
  - Voldemort, Riak, BigCouch
  - MongoDB, Redis, CouchDB, Neo4J
- Associated use case:

COMPUTERWORL

• All of the above

 $\ensuremath{\mathbb{C}}$  2011 by The 451 Group. All rights reserved



# Necessity

- BigTable: Google
- Dynamo: Amazon
- Cassandra: Facebook
- HBase: Powerset
- Voldemort: LinkedIn
- Hypertable: Zvents
- Neo4j: Windh Technologies
  - Yahoo: Apache Hadoop and Apache HBase
  - Digg: Apache Cassandra
  - Twitter: Apache Cassandra, Apache Hadoop and FlockDB





## **Introducing Total Data**

- Defined by The 451 Group to describe new approaches to data management
- Reflects the changing data management landscape as pragmatic choices are made about data storage and analysis techniques
- Processing any data applicable to analytic query
  - in database, data warehouse, or Hadoop, or archive
  - structured or unstructured, relational or non-relational
  - on-premise or in the cloud
- Inspired by 'Total Football'





#### **Total Football**

 "You make space, you come into space. And if the ball doesn't come, you leave this space and another player will come into it."

Bernadus Hulshoff, Ajax 1966-77





COMPUTERWORLD

#### Data management – in theory

- The application is the primary source of data
- The relational database is sacrosanct
- The enterprise data warehouse is the single source of the truth (or is supposed to be)
- Data archived to tape
- Infrastructure primarily exists to support the data/application layer



**COMPUTERWORLD** 

#### Data management – in practice

20.

- The relational database is sacrosanct
- Distributed data layer to meet the scalability and performance demands
- New opportunities for real-time BI
- Polyglot persistence use the most appropriate data storage for the application





#### **Databases in the cloud**

- 2008: Operational database moved to the cloud: EnterpriseDB, MySQL, Oracle, SQL Server, DB2 (09)
- 2009: Analytic databases moved to the cloud: Vertica (08), Aster Data, Netezza, RainStor, Greenplum, Teradata
- Some dev and test, some early adopters, but limited take-up
- The software, and licensing, is not designed to be elastic
- Data security, management concerns
- A true "cloud database" must be designed to take advantage of elastic scaling and multi-tenancy



M all of

### **Cloud databases**

- 2008: Introduction of SimpleDB, Microsoft SDS, Google App Engine
- 2009: Amazon RDS, Microsoft SQL Azure, Rackspace FathomDB
- 2010: NoSQL database vendors: 10gen, DataStax, Hypertable, Basho, Couchbase, Neo Technology,
- 2010: Google High Replication Datastore, NewSQL projects.
- 2011: More from Amazon, Google, Microsoft, Rackspace, NewSQL vendors, distributed data grid providers GigaSpaces CEAP
- 2012: True cloud databases emerge. PaaS emerges as a potential battleground for NoSQL, NewSQL and DDG providers



#### Data location, data location, data location

- Not the end of the EDW, but EDW is one of many sources of BI, rather than the only source of BI
- Choose the right storage technology: software and hardware
  - EDW, Hadoop or archive
  - On-premise or on the cloud
  - Memory, disk or SSD



#### Data location, data location, data location

- Understand the requirements:
  - Value and temperature of the data
  - Ensure data can be queried using existing tools/skills
  - Cost
- The issue of data location becomes paramount
  - Avoid data movement and duplication retain governance
  - Virtual data marts and data cloud
  - Data virtualization to provide access to multiple data sources



#### **Relevant reports**

#### Total Data

COMPUTERWORLD

- Explaining the total data management approach to dealing with the impact of big data on the data management landscape
- Coming late 2011
- Including the growing Hadoop ecosystem and real-time
- sales@the451group.com



M the st

# Conclusion

- The key factors driving the adoption of alternative data management technologies are scalability, performance, relaxed consistency, agility, intricacy and necessity (SPRAIN).
- NoSQL projects were developed in response to the failure of existing suppliers to meet the performance, scalability and flexibility needs of large-scale data processing, particularly for Web and cloud computing applications.
- NewSQL and data-grid products have emerged to fulfill a similar requirement among enterprises, a sector that is also now being targeted by NoSQL vendors.

M Cho and

COMPUTERWORL

# Conclusion

- The database market remains dominated by relational DBs and incumbent industry giants, but the rise of NoSQL and NewSQL has in part been driven by the failure of these products to address emerging data management requirements.
- For the most part these database alternatives are not designed to directly replace existing database products, but to offer purposebuilt alternatives for workloads that are unsuited to general-purpose relational databases.
- However, PaaS and the distributed data layer that will power future
   PaaS is emerging as a potential battleground for NoSQL, NewSQL and
   DDG providers.

© 2011 by The 451 Group. All rights reserved

