



## **Analytic Platforms: Beyond the Traditional Data Warehouse**

**By Merv Adrian and Colin White**

**BeyeNETWORK Custom Research Report  
Prepared for Aster Data**

## Executive Summary

The once staid and settled database market has been disrupted by an upwelling of new entrants targeting use cases that have nothing to do with transaction processing. Focused on making more sophisticated, real-time business analysis available to more simultaneous users on larger, richer sets of data, these analytic database management system (ADBMS) players have sought to upend the notion that one database is sufficient for all storage and usage of corporate information. They have evangelized and successfully introduced the **analytic platform** and proven its value.

A dozen or more new products—the majority introduced after 2005—have been launched to join the pioneering analytics-specific offerings, each of which boasts thousands of installations. Collectively, the newcomers successfully placed an additional thousand instances by the end of the decade, making it clear that the analytic platform has tapped into a significant market need. They have added hundreds of millions of dollars per year to the billions already being spent with the early entrants—and taken share from incumbent “classic data warehouse relational database management system” products.

Analytic platforms provide two key functions: they manage stored data and execute analytic programs against it. We describe them as follows:

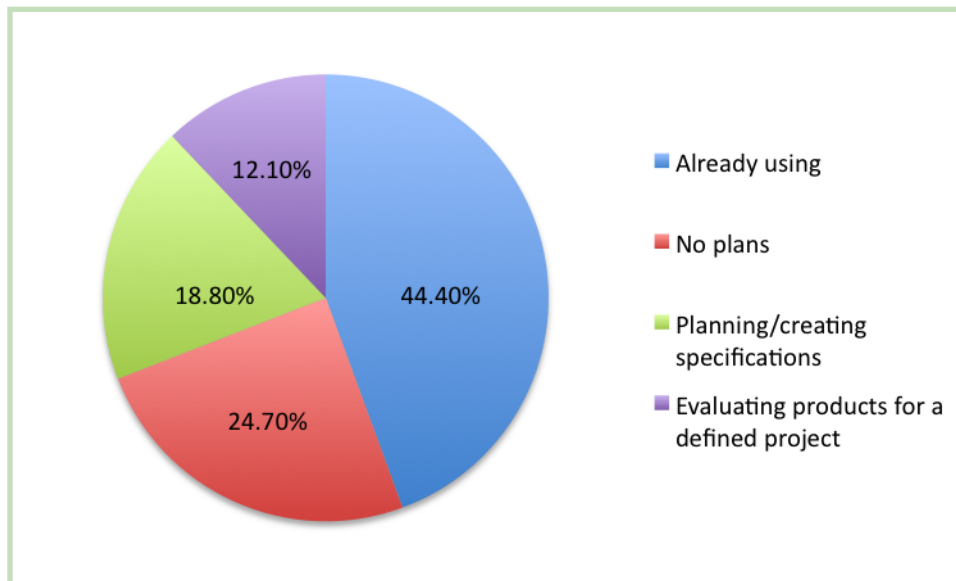
*An analytic platform is an integrated and complete solution for managing data and generating business analytics from that data, which offers price/performance and time to value superior to non-specialized offerings. This solution may be delivered as an appliance (software-only, packaged hardware and software, virtual image), and/or in a cloud-based software-as-a-service (SaaS) form.”*

Some survey respondents, when confronted with this definition, disagreed with it—they consider the “platform” to be the tools they use to perform the analysis. This may be a legacy of client-server days, when analysis was performed outside the database on “rich client” software on desktops. But the increasing requirement for the ADBMS to power the analysis is upending this thinking, and most agreed with our description. We found:

- **The pace of adoption is strong and accelerating.** In 2009, thousands of analytic platforms were sold. And 10 or more players with growing sales are competing for an increasing number of use cases, worldwide, in many industries.
- **The promises being made are being met.** Adopters of analytic platforms report that they tested difficult problems in proof-of-concept (POC) exercises, and the selected products were equal to the tasks—beating their incumbent DBMSs.
- **The right selection process is essential.** Successful POCs require an understanding of the likely analytical workloads—data types and volumes, the nature of the analysis, and the numbers of users likely to be on the system. And real tests separate winners from losers: often, some candidates can’t get it done at all.

## Introduction

We conducted an online survey of several hundred professionals worldwide, who shared their experiences and opinions with us. Survey results are shown at the end of this report; we include some highlights throughout. Only 25% of those surveyed said they have no plans for an analytic platform. 44% said they are already using one (see Figure 1).



**Figure 1: Are You Using or Planning to Use an Analytic Platform?**

We also conducted interviews with 8 analytic platform vendors, all of whom are targeting this market, and with a nominated customer from each. The interviewees are quite different from the overall survey population. While our survey showed organizations using database management system (DBMS) products for analytic platforms in proportions that mirrored overall market shares, our interviewees come from the leading edge of the disruptive analytic platform phenomenon. They work for organizations that continue to use classic relational database management systems (RDBMSs) for many applications, including some of the business analytics being targeted by the vendors of analytic platforms, but have opted to use specialty platforms for a variety of reasons.

What we learned was profound; businesses, more and more driven by their need for analytic processing of enormous amounts of data, are responding to the emergence of a class of DBMS specialized for analytics, recently introduced to the market in most cases. A thousand sales of these products in just a few years, generating billions of dollars in revenue, herald the arrival of the analytic platform as a category to be watched closely. It solves important problems, and customers are deriving enormous value from it, creating new classes of business applications and driving top-line growth.

Our interviewees were unanimous: their money was well spent, and their existing classic RDBMS offerings fell short. By contrast, only 21.4% of 168 survey respondents, many still using classic RDBMS products for their analytic platforms, pronounced themselves fully satisfied with their analytic platform projects. While we did not ask for their reasons for this dissatisfaction, some can be derived from the “issues that led you to add an analytic platform” data: the need for complex analyses, query performance, and on-demand capacity topped the list. These issues are mirrored in the case study interviewees.

This report examines the analytic platform, the business needs it meets, the technologies that drive it, and the uses analytic platforms are being put to. It concludes with some guidance on making the right choices and getting started with the products of choice.

## The Business Case for Analytic Platforms

### What is an Analytic Platform?

Informally: the analytic platform is a response to the inadequacy of classic RDBMSs for new, more sophisticated demands for business analytics. It combines the tools for creating analyses with an engine to execute them, a DBMS to keep and manage them for ongoing use, and mechanisms for acquiring and preparing data that is not already stored. In this report, we focus on the DBMS component of the platform. As noted below, separate providers also offer data sourcing and integration and tools for analytics surrounding the DBMS; these will interact with the DBMS itself and often depend on it for execution.

### Why Do We Need Analytic Platforms?

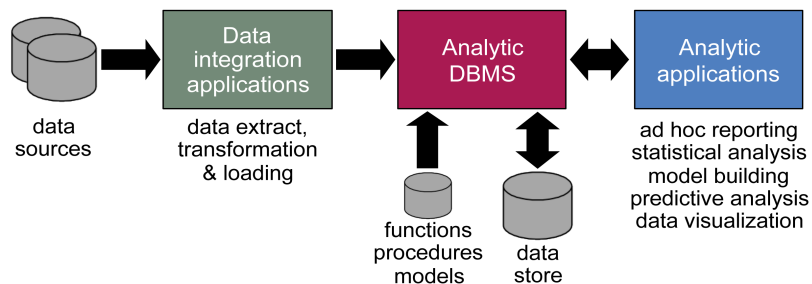
A brief history demonstrates how we got here over several decades. The earliest computing used a simple paradigm for simple analytic processing: business data created by transaction, manufacturing, or other processes was stored in files. Specialists wrote programs to run against them, generating management reports about the state of the business. But routine, multiple, simultaneous use of the data—transactional and reporting—quickly became the expectation.

DBMSs emerged: persistent data stores for many kinds of batch programs to run against—to add, update and delete, and report on data. Online computing made it possible to do these things in real time, and to do them at the same time—multiuser multiprogramming. The client-server era shifted things to a two-or-more-tier model, in which the analytic processing was done on data extracted to a different platform, supporting one or many users working with local copies of the data that might themselves be saved or might go away when the session was done. But this created uncoordinated, redundant, and sometimes conflicting versions of the data.

The data warehouse was envisioned as a central data store where access, definitions, governance, policy, and currency could be centrally managed. Diverse data sources were harvested and data was copied in, separating analytics and reporting from other business processing. Over time, satellite data marts for specific subject areas or user populations or both emerged—along with rising budget authority in business units who desired autonomy. “In front” of these systems, data extraction and transformation products managed feeding the data in; “behind” them, analytic tools for ad hoc reporting, statistical analysis, model building, predictive analysis, data visualization, etc. were created for business users, programmers, and non-programmers alike, to use (see Figure 2). But the DBMS product in the middle of all this was usually the same one in use for everything else.<sup>1</sup>

---

<sup>1</sup> Teradata, 4GLs like FOCUS and SAS and other products in the 80s were positioned as “storage plus analysis” vehicles for large volumes of data. But most buyers considered their standard, classic RDBMS as the default.



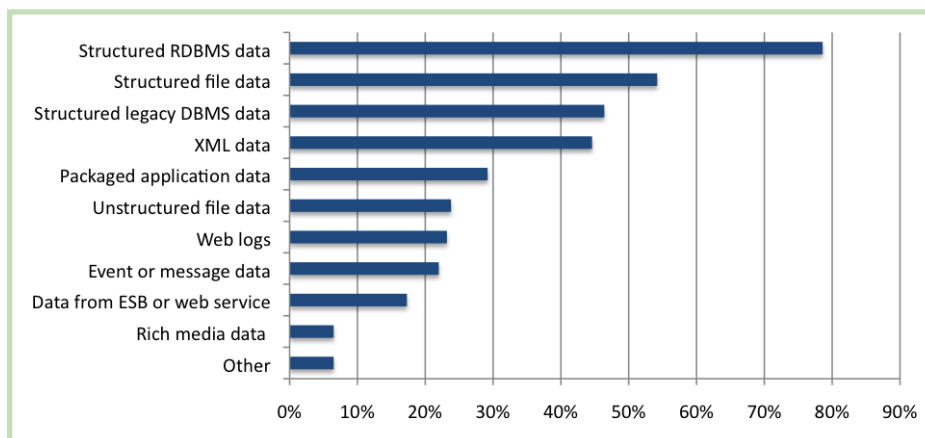
**Figure 2: Components of an Analytic Platform**

Requirements continued to become more difficult to meet. Online analytic processing (OLAP) added multidimensional capability, the ANSI/ISO SQL standards steadily added more power to the language used in databases, and the TPC-H benchmark was created to measure analytic performance. The benchmark made it clear that DBMSs were coming up short; new approaches were needed, and new vendors emerged to meet them, creating new products that succeeded where the incumbents could not. The forces driving the need for change are largely the same and drove the design of the newcomers who joined the pioneering offerings from Sybase and Teradata.

### Data Growth and New Types of Data

The largest data warehouses are now measured in petabytes. Terabytes are not at all unusual, and it's routinely reported that the largest are growing at an increasing rate—tripling in size every 2 years. Sixteen percent of the analytic platforms reported in our survey were managing more than 10 terabytes after loading, tuning, enhancing, and compressing it, and 64.3% said that support for more than 100 terabytes was very or somewhat important in their planning or acquisition.

And while data warehouse data is mostly structured, a significant amount of other corporate data is not. Unstructured text data, weblogs, scientific feeds, photographs, videos, sound files are all potential sources for analysis. Our survey respondents are feeding their analytic platforms with a variety of these new data types: 44.6% are using XML data, 23.8% unstructured file data, 23.2% weblogs, and others (see Figure 3). But the languages, analytic products, and storage mechanisms that have been the everyday toolkit for business analysts were not designed for these new forms of information and often are not well-equipped to work with them.



**Figure 3: What Data Sources are Used in Your Analytic Platforms?**

Analytic platforms are designed to manage large data volumes, sophisticated analytics, and newer data types. They use modern storage paradigms that allow retrieval by columns more efficiently and encode data for better compression. Some use “smart storage” to do some of the work at the storage layer to free up the processor for the heavy analytic lifting. They lash together many commodity processors, with larger memory spaces. They connect processors with one another and with data storage across faster networks to scale processing and storage in sync. They are designed to handle new types of data, even user-defined types that may have specialized code associated with them to make sense of the unfamiliar content. Classic RDBMS products were not built with these innovations in mind and are not always easy to update to leverage these new opportunities.

### **Advanced Analytics**

Simple reporting, spreadsheets, and even fairly sophisticated drill-down analysis have become commonplace expectations and are not considered “advanced.” While the term is frequently debated, it’s clear that even “simple” analysis is advanced when it needs to be performed on a massive scale. Even a simple comparison of profitability across 2 days’ trade activity for the top 10 traders each day, for example, is a performance challenge for many systems when run against today’s extraordinary volumes of data while other activities run on the same system.

But increasingly, the nature of the analysis itself is more “advanced.” Sophisticated statistical work is becoming commonplace for market basket analysis in retail, behavioral analysis in clickstreams for websites, or risk analysis for trading desks. Building predictive models and running them against real-time data is a frequent use case. Some firms require geographic visualizations of data, often against variable shapes such as sales territories or watersheds that are not easily computed. Such ambitions used to be left to the largest firms with highly sophisticated programmers and expensive hardware in their data centers. No longer; savvy business leaders, even in mid-size firms, expect the same from their teams today. And they are doing it outside their classic RDBMS; in our survey, 53% of 223 respondents said they perform business analysis on data not contained within an RDBMS. Nearly two-thirds of them were using hand-coded programs as opposed to packaged tools.

Analytic platforms address the specialization implicit in handling analytic workloads. They retrieve and manipulate large sets, using the right subsets of the fields in individual records. They support large memory spaces for this processing, dramatically improved I/O times to get the data there from storage, and support not only advanced SQL’s capabilities but also user-defined functions (UDFs) and programming languages that analysts and statisticians often use instead of SQL. And they leverage new paradigms like MapReduce programs, which may run over external files or against data imported from those sources.

### **Scalability and Performance**

Data scalability is only one dimension—the other is multiuser performance. It has long been a goal of business intelligence (BI) thinkers and planners to involve more users in the corporate analysis of performance. In the client-server era this was often handled by putting tools on their desktops and moving data to them, creating coordination problems as computational models were duplicated. Unsynchronized, often contradictory analysis resulted.

Centralizing the key metrics and algorithms, and making them consumable by more employees and partners who can collaborate around their work, are key challenges. Our survey users expect high volumes of simultaneous usage—32.1% say they need to support more than 100 concurrent users.

Analytic platforms are designed to leverage higher bandwidth connections across a fabric of processors. They utilize modern “container” constructs in memory, used to protect and coordinate multiple processes running in massively parallel scale-out architectures with more processors. They use inexpensive hardware that can be added without taking systems down, so as demands scale, so can processing power. They are designed to cooperate with virtualization layers in modern environments that permit the elastic setup and teardown of “sandboxes” where new analyses and ideas can be tested. All of these capabilities permit analytic platforms to raise the performance profile.

### **Cost and Ease of Operation**

As data volumes, analytic complexity, and the numbers of users all grow, so does cost. Even “commoditized” hardware costs millions of dollars; capital costs expand with data, power, and populations. Power, cooling, space, and backup/recovery for all of it add more expense. Moreover, additional disks and more processors mean more management. Policies across multiple classes of users, security management, and the need to manage environments that cannot be taken down for maintenance all create their own demands and costs.

The number of moving parts in these systems creates its own added challenge: the difficulty of “standing the system up” in the first place becomes an exercise in coordinating software versions, device drivers, and operating systems. Each piece of a complex stack of software is frequently updated by its supplier—and one piece’s fix breaks another piece. Systems management skills become expensive, and budget is consumed merely “keeping the lights on.”

Analytic platforms offer multiple deployment options that can reduce many of these costs. As they generally move to commodity hardware, some of the pricing premium in older proprietary systems is eroded. The replaceable form factor of massively parallel processing (MPP) systems makes scaling smoother and more granular. It is simpler to add blades with processor, memory, and storage that snap into racks and can be bought as needed. Open source software used in many stacks lowers licensing costs.

Appliances—pre-integrated, preconfigured collections of hardware and software or bundles of multiple software parts that may be installed on any commodity hardware system—offer a way to reduce setup cost. They are increasingly maintained and updated by their suppliers in a way that is designed to ensure that changes don’t “break things.”

Finally, moving the analytic platform off premises in one fashion or another provides the maximum reduction in cost of ownership and operation. Several vendors will host the system and the data as a dedicated facility. Some will make it available “in the cloud” in a multi-tenant fashion, where tools are shared but data is stored and managed for individual customers. They may take over the process of importing the data from its source systems, such as retail or online gaming systems, and provide the data integration as well as the storage and analytics.

Recall our formal definition:

*An analytic platform is an integrated and complete solution for managing data and generating business analytics from that data, which offers price/performance and time to value superior to non-specialized offerings. This solution may be delivered as an appliance (software-only, packaged hardware and software, virtual image), and/or in a cloud-based SaaS form.*

In this report, we consider DBMS offerings that form the heart of the analytic platform.

## Types of Analytic Platforms

For the past few decades, RDBMS products have formed the data management underpinnings of a wide range of both transaction and analytic IT applications. Products that target analytic processing can be thought of as ADBMSs. Some ADBMS products support SQL and the relational model, while others offer alternative languages and data models.

There are numerous features and functions that differentiate ADBMSs from one another, but for the purposes of simply describing the players, they may be classified in several key dimensions:

- **Use of proprietary hardware:** Some vendors create their own specialized hardware to optimize processing. Others run on any standard hardware.
- **Hardware sharing model for processing and data:** Increasingly, ADBMS vendors support MPP architectures which distribute processing across many blades using chips with multiple cores and significant amounts of dedicated on-board memory. These may have dedicated storage in a shared-nothing environment or may be connected to shared storage such as a storage area network (SAN).
- **Storage format and “smart data management:”** Many ADBMSs are using columnar storage, which dramatically improves the performance of disk I/O for certain read operations. Some support both row and column format in one hybrid form or another. Some also add intelligence at the storage layer to pre-process some retrieval operations. All use a variety of encoding, compression and distribution strategies.
- **SQL support:** Support for “standard” SQL tends to depend on which standard you mean; no vendor supports all of the SQL languages. The absence of some specific features, like correlated subqueries or joins across tables on separate nodes, can be a serious performance problem, preventing some queries from running adequately or at all.
- **NoSQL too.** Recently, a number of offerings have emerged for analyzing specific data types such as documents, unstructured text, and other content not typically stored inside RDBMSs. These are often collectively referred to as NoSQL solutions. Some actually store data while others, such as MapReduce, may operate on files stored in a file system like the open source Apache Hadoop Distributed File System. These offerings are relatively specialized at this time, but can be very effective. Many are adding more features that provide data import, SQL query, and other RDBMS-like functionality.
- **Programming extensibility:** ADBMS engines offer varying degrees of support for the installation of functions, usually described as UDFs, which offer callable computations and data manipulations that are difficult to reproduce with standard SQL. Some offer libraries of such functions themselves and with partners, and some of these take advantage of system parallelism for performance improvement.
- **Deployment models.** ADBMSs may be delivered as an appliance: a complete package of hardware and software; software-only products may be deployed on premises on commodity hardware, hosted off premises, or even in public clouds such as Amazon’s EC2.

## Hardware Directions

Things are changing fast. Several key elements of the hardware mix are undergoing enormous change, with profound implications for system design and its impact on analytic performance.

**Memory is the new disk; disk is the new tape.** Reading and analyzing data is made much easier when the data all fits in memory; disk I/O problems, the management of buffers, writing out to disk when new data needs to be brought in—all of these become less of a performance challenge. Memory prices continue to drop and the use of solid state disks (SSDs) and flash memory are rewriting the rules. The first all-memory systems are already appearing, and more will come.

**More cores, more threads, yield more processing power.** The addition of more cores (and processing threads) to chips has similar implications. As software smart enough to break up and distribute the work (parallelization) is given more threads to work with, performance can scale simply with the addition of more standard blades to a system. In MPP systems where storage is dedicated to the processor, this scalability extends not just to power or number of users but also to data volume.

**Infiniband and other network interconnects drive speed.** The speed of interconnects can be an enormous bottleneck for system performance. Moving data around inside large systems or from one system to another becomes more difficult with larger volumes. Infiniband's raw speed and ability to provide parallel data movement will be a key asset for vendors that utilize it.

## Message from the Market: It's Time

Markets change rapidly, but the effects are often not felt for years. The value of already installed software in most categories is several orders of magnitude larger than the spending on it in any given year or two. Maintenance and support costs for installed software dwarfs new spending. But at the leading edge, players and industry analysts are dazzled by new products and new sales.

Analytic platforms are no exception to this. From the mid-1990s to the mid-2000s, Sybase and Teradata were largely alone in the specialty analytic database market. By 2010, they had some 6,000 installations of their products between them. A dozen or so newer vendors, emerging throughout the last decade, added another thousand or so. The several hundred million dollars spent with these newcomers represented the most significant spending shift in database systems in decades.

But in context, these numbers are hardly a blip on the radar. There are hundreds of thousands of DBMSs installed; so-called data warehouse DBMS sales are estimated at \$7 billion per year. The ADBMS is in the hands of early adopters, not mainstream customers—even when they are being used by the world's largest enterprises, their use is confined to a business unit, a division, or a team of specialists. Leaving aside Teradata and Sybase, ADBMS vendors collectively generate a few hundred million dollars annually—less than 5% of the data warehouse DBMS market. Small wonder, then, that our survey respondents told us that they typically begin their search for a platform with their incumbent DBMS vendor.

We learned in our interviews that those adopting analytic platforms are agents of change. They are creating new value, new business opportunities, and new customer opportunities. From a competitive point of view, organizations that have not yet assessed ways to leverage these platforms are already

behind. That's the bad news. The good news? One of the key findings of this report is that if you know your problem, you can start fast. And get value fast. At lower cost than you may have thought possible.

## Techniques and Technologies

In this section, we review some of the key techniques and technologies offered by analytic platforms, and offer some suggestions about things to consider when evaluating these solutions.

### ADBMS versus a General Purpose RDBMS

An analytic platform consists of three main software components: the data integration software for transforming and loading source data into the platform's database, the database management software for managing that data, and the analytic tools and applications that analyze the data and deliver analytics to users. In a traditional data warehousing environment, these three components are purchased separately and integrated by the customer. A key difference with an analytic platform is that the vendor does the integration and delivers a single package to the customer.

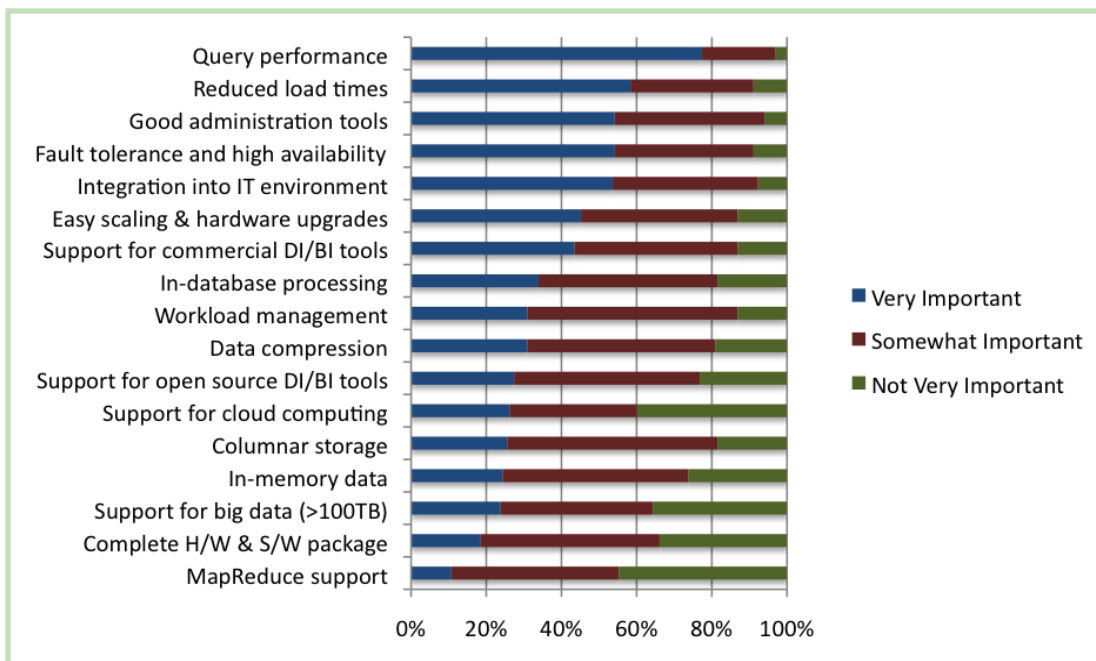
At present, most analytic platform database management is done by RDBMSs. For the past few decades, RDBMS products have formed the data management underpinnings of a wide range of both transaction and analytic IT applications. Given the trend by many companies toward extreme processing at both the transaction and analytic ends of the application processing spectrum, it is becoming more difficult for a general purpose or *classic* RDBMS to support the increasing number of different uses cases and workloads that exist in organizations.

The broadening application processing spectrum is leading to vendors developing database management software that focuses on a narrower subset of that spectrum. In this report, products that target analytic processing are described as ADBMSs.

Even within the ADBMS segment, the ability of any given product to support a specific use case or workload varies. The challenge in selecting an ADBMS is to match the workload to the product. This is especially true in the case of extreme processing and also in business environments with constantly changing requirements. Often the only solution is to run a POC evaluation using real user workloads.

In our study, we focused primarily on ADBMS solutions that support the relational model and SQL. However, a brief discussion on using a non-relational, or NoSQL, approach is also included.

In our survey and customer interviews we asked people about key technology requirements for an analytic platform. Our objective was to determine the characteristics and features of an analytic platform that were most important to organizations. The survey results are shown in Figure 4. Features rated as *very important* by the majority of respondents were: query performance (77%), reduced load times (58%), good administration tools (54%), fault tolerance and high availability (54%), and integration into the existing IT environment (54%). Other features that received high scores were easy scaling and hardware upgrades (45%), support for commercial data integration and BI tools (44%), and in-database processing (34%).



**Figure 4: What Features are Important for Your Analytic Platform?**

There were no major surprises in these scores except that the score for *in-database processing* was higher than expected. This demonstrates that organizations are beginning to appreciate the benefits of exploiting the power of a parallel database engine to run certain performance critical components of a data integration or analytic application.

Scores that were lower than expected were: support for open source data integration and BI tools (27%), columnar data storage (26%), and a complete hardware and software-packaged solution (18%). The first two results may reflect the limited experience of organizations in using these technologies, while the third score demonstrates that respondents often prefer to have the flexibility to choose their own hardware.

Of the 8 customers interviewed for the report, 7 were doing *extreme processing* involving significant amounts of detailed data and intensive SQL processing. All 7 stated that query and load performance coupled with easy scaling were the main product selection criteria. Most of these customers also required high availability.

### ADBMS Application Development Considerations

When RDBMS technology and SQL were introduced in the early 1980s, the big leap forward was separating the user, or logical, view of data from the way it is physically stored and managed. An RDBMS optimizer handles the mapping of SQL requests to the physical storage layer. This explains why the quality of a product's optimizer can play a big role in performance. Even today, this physical data independence remains largely unique to relational technology.

From a development perspective, the factors to consider when selecting an analytic platform and its underlying ADBMS are: its SQL functionality, the programming languages supported, the quality of the relational optimizer, and the physical storage options provided. Some of these factors are of more

concern to applications developers than the users of interactive analytic tools. For these latter users, the main consideration is whether the SQL support provided by the ADBMS is sufficient to allow the analytic tool to operate efficiently.

As already noted, most of the customers interviewed for this report were using extreme processing, and in all these cases, a certain percentage of the end users were creating their own ad hoc SQL queries. These queries were often very sophisticated, and the analytic platform's SQL support was a very important selection criterion for these customers. Several customers commented that some of the products they evaluated during product selection had inadequate SQL functionality. Also, with certain products, the physical layout of the database imposes restrictions on the SQL that can be used, which of course is contrary to one of the main tenets of the relational model.

One major area of difference between vendors is their support for SQL built-in functions (scalar, aggregate, string, statistical, etc.), UDFs, stored procedures, and other types of in-database processing such as MapReduce, predictive models, etc. The ability to *push* analytic functions and processing into the ADBMS will usually boost performance and make complex analyses possible from users who have the expertise to use such functions, but not the skills to program them. For many of the customers we interviewed, in-database processing was an important feature when choosing a product. The use of such processing, however, can limit application portability between different ADBMS products because of implementation differences.

It is important to note that just because an ADBMS product supports a particular type of in-database processing, it does not necessarily mean this processing is done in parallel. Some of the processing functions may be run in parallel, while others may not. All of them provide more rapid implementation, but the parallelized ones offer superior performance. As an example, not all products support the ability to store and run multiple copies of the same stored procedure on multiple nodes of the configuration.

### **ADBMS Data Storage Options**

ADBMS software supports a wide variety of different data storage options. Examples include: partitioning, indexing, hashing, row-based storage, column-based storage, data compression, in-memory data, etc. Also, some products support a shared-disk architecture, while others use a shared-nothing approach. These options can have a big impact on performance, scalability, and data storage requirements. They also cause considerable discussion between database experts as to which option is the best to use. The current debate about row-based versus column-based storage is a good example here. Often these debates are pointless because different products implement these features in different ways, which makes comparison difficult.

In an ideal world, an ADBMS would support all these various options and allow developers to choose the most appropriate one to use for any given analytic workload. ADBMS products, however, vary in their capabilities. Of course, providing too many alternatives adds complexity to the product, to application deployment, and to database administration. A product could automatically select or recommend the best option, and some products are beginning to support this. In general, however, this type of feature is difficult to implement successfully given the complexity of today's analytic workloads.

The physical storage options supported by an ADBMS product should be completely transparent to the user's view of the data, i.e., the user should not be forced to code SQL queries to suit the way the data

is physically stored. Realistically, in the case of extreme processing, some tuning of SQL queries and the building of indexes and aggregates common in classic RDBMSs may still be necessary to obtain the best performance. This was certainly the case for several of the customers interviewed for the report.

Another option of course is go with a product that provides very little in the way of tuning options and instead employ a brute-force approach of simply installing more hardware to satisfy performance needs. The theory is that hardware today is cheap compared to development and administration costs. This is often the approach used in NoSQL products.

### **The Role of MapReduce and NoSQL Approaches**

No single database model or technology can satisfy the needs of every organization or workload. Despite its success and universal adoption, this is also true for RDBMS implementations. This is why some organizations develop their own tailored solutions to address certain specific application needs.

Google is a good example. Like many other Internet-based organizations, Google has to manage and process massive amounts of data every day. A high percentage of this data is not well structured and does not easily lend itself to being managed or processed by a RDBMS. To solve this problem Google developed its own technology. One important component of this technology is a programming model known as MapReduce.

A landmark paper<sup>2</sup> on MapReduce by Jeffrey Dean and Sanjay Ghemawat of Google states that:

*“MapReduce is a programming model and an associated implementation for processing and generating large data sets .... Programs written in this functional style are automatically parallelized and executed on a large cluster of commodity machines. The runtime system takes care of the details of partitioning the input data, scheduling the program’s execution across a set of machines, handling machine failures, and managing the required inter-machine communication. This allows programmers without any experience with parallel and distributed systems to easily utilize the resources of a large distributed system.”*

MapReduce programs manipulate data records that are formatted as *key/value* pairs. The records are produced from source data by the map program. The *value* field of a data record can contain any type of arbitrary data. Google uses this approach to index large volumes of unstructured data. Note that MapReduce is not a new concept—it is based on the list processing capabilities in programming languages such as LISP (LISt Processing).

The MapReduce programming model has now been implemented in several file and database management systems. Google has integrated it into its BigTable system, which is a proprietary DBMS that uses the Google File System (GFS). It is also a component of the Apache open source Hadoop project, which enables a high-scalable distributed computing system. In the case of Hadoop, MapReduce is deployed on the Hadoop Distributed File System (HDFS).

The MapReduce programming model has also been implemented in a number of ADBMS products. Several of the sponsors of this report provide this capability. This hybrid approach combines the advantages of the MapReduce programming model with the power and integrity of a parallel database engine.

---

<sup>2</sup> <http://labs.google.com/papers/mapreduce.html>

The advent of MapReduce has led to the development of a wide variety of solutions that offer alternatives to RDBMS technology. This group of solutions is often referred to as the *NoSQL movement*. These solutions include not only products that support MapReduce processing, but also document and XML data, graph data, etc. Examples of software here include Amazon Dynamo storage system, Apache Cassandra (originally developed by Facebook) and CouchDB projects, MarkLogic, and MongoDB.

The availability of NoSQL software has led to a heated debate about the pros and cons of these solutions vis-à-vis RDBMSs. The NoSQL advocates say that NoSQL solutions are superior to RDBMSs and will ultimately replace them, whereas the RDBMS camp say the NoSQL software lacks integrity and reliability.

The NoSQL debate is reminiscent of the object-relational database wars of the 1980s. The reasons behind them are similar. Programmers prefer lower-level programmatic approaches to accessing and manipulating data, whereas non-programmers prefer higher-level declarative languages such as SQL. The inclusion of MapReduce in ADBMS products offers some of the best of both worlds.

One issue with NoSQL technology is that some software organizations are reinventing the wheel by trying to extend NoSQL software with features that RDBMS vendors have spent many years refining and optimizing. In some cases NoSQL software developers are even adding SQL support. A better solution is to recognize that both technologies have their benefits and to focus instead on making the two coexist together in a hybrid environment.

Many ADBMS and NoSQL solution providers agree that enabling a hybrid environment is what most customers want, and are building connectors between the two technologies. Maybe this is why the website *nosql-database.org* prefers the pragmatic term *Not only SQL* to NoSQL.

MapReduce is particularly attractive for the batch processing of large files of textual data. Seven percent of our survey respondents were using MapReduce with Hadoop. One of the customers interviewed for our study was using a hybrid environment where Hadoop and MapReduce were used for processing textual data, and subsets of this data were then brought into the analytic environment using a software bridge from Hadoop to the ADBMS.

### Administration and Support Considerations

Good administration capabilities rated high in our survey results (54% rated it as very important) and customer interviews. Several of the customers interviewed also said that simple administration was an important product selection criterion because they didn't want to employ "an army of database administrators." Easy administration was particularly important when designing databases and storage structures, and when adding new hardware.

Several of the customers interviewed also noted that as workloads increased in volume and became more mixed in nature, the workload management capabilities of the ADBMS became more important. Some said they wished they had done a better job of testing out mixed workloads in POC trials.

All of the interviewed customers were happy with the support they received and the working relationship they had with their vendors. Several also commented that the vendor was usually very receptive to adding new features to the analytic platform to meet their needs.

## Deployment Models

The deployment options offered by analytic platform vendors vary. Some vendors provide a complete package of hardware and software, while others deliver an integrated pack of software and then let customers deploy it on their own in-house commodity hardware. Some vendors also offer virtual software images that are especially useful during for building and testing prototype applications.

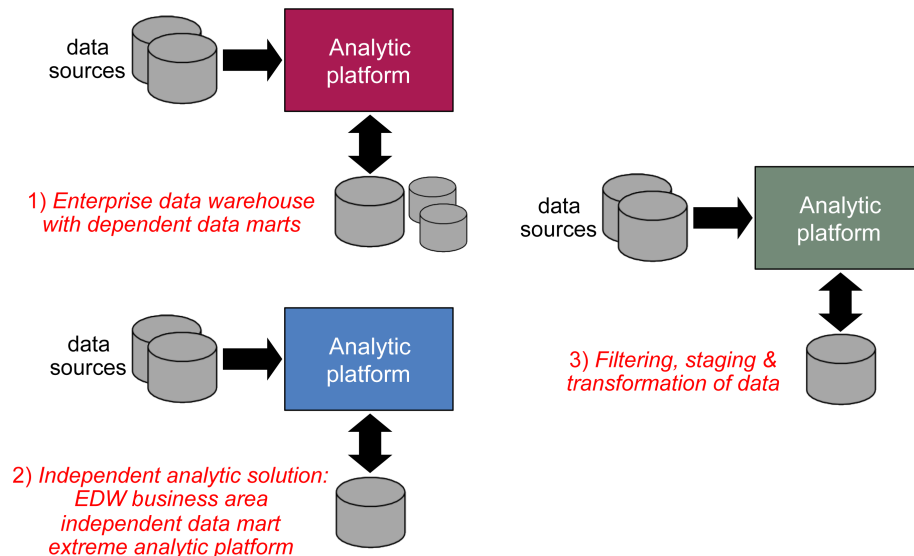
One direction of the analytic platform vendors is to provide cloud-based offerings for deployment in the either the vendor's or a third-party data center or for use on an in-house private cloud. In some cases, the vendor may also install and support a private cloud analytic platform on behalf of the customer.

Ideally, a vendor should support a variety of different deployment options for its analytic platform. This gives customers the flexibility to use the most appropriate environment for any given situation. The customer may opt, for example, to develop and test an application in a public cloud and then deploy the application in house. Other customers may wish to use a hybrid environment where some applications are run in house, while others are deployed in a public cloud depending on performance, cost and data security needs.

## Use Cases

Based on prior experience, the survey results and customer interviews from our research study, we can identify three dominant use cases for an analytic platform (see Figure 5):

1. Deploying an enterprise data warehousing environment that supports multiple business areas and enables both intra- and inter-business area analytics.
2. Enabling an independent analytic solution that produces analytics for an individual business area or to satisfy a specific business need
3. Facilitating the filtering, staging, and transforming of multiple data sources and types of data for use in analytic processing



**Figure 5: Analytic Platform Use Cases**

Before looking at each of these use cases in detail, it is important to comment about the survey results and customer interviews used in this section of the report.

The organizations and users surveyed represent a wide spectrum of industries, data warehousing environments, and technology maturity. The customers interviewed for the report, on the other hand, were recommended by each of the report sponsors, and were, in many cases, developing analytic solutions where it was not practical to maintain the data in a traditional data warehousing environment.

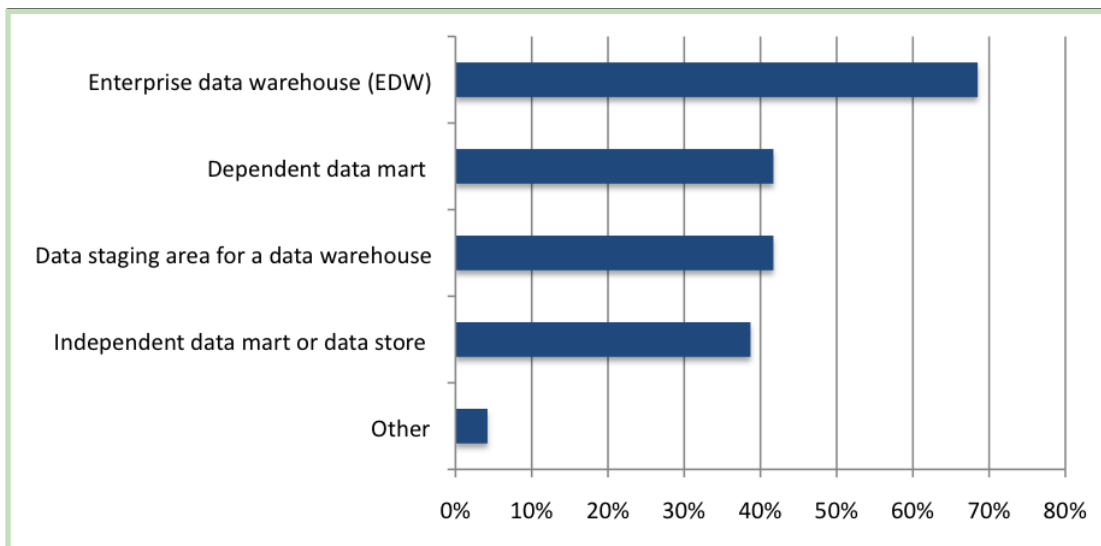
The results and opinions from the two groups therefore sometimes differ. The survey audience results reflect the ongoing evolution of the traditional data warehousing environment, whereas the opinions of the interviewed customers demonstrate the disruptive forces taking place in the industry that enable completely new types of analytic application to be developed.

### Use Case 1: Enterprise Data Warehousing

This use case is well established and represents what can be considered to be the *traditional* data warehousing approach. The environment consists of a central enterprise data warehouse (EDW) with one or more virtual or dependent data marts. The data in the EDW and data marts has been cleansed and transformed to conform to IT designed data models and may be kept almost indefinitely for historical reporting and data analysis purposes by multiple business areas.

The survey results (Figure 6) showed that 68% of survey respondents were using an analytic platform for deploying an EDW, while 42% were using the analytic platform for a dependent data mart containing data extracted from an EDW.

Of the 8 customers interviewed for this report, only one was using an analytic platform for enterprise data warehousing. For this customer, reducing software costs was the main reason for moving to an analytic platform from a classic RDBMS product (i.e., a database system that is used for both transaction and analytic processing).



**Figure 6: What Use Cases are Being Deployed on Your Analytic Platform?**

## Use Case 2: Independent Analytic Solution

In this use case, the data being analyzed is maintained outside of the EDW environment. Some 39% of survey respondents were using an analytic platform for this use case. There are three main reasons why an organization may choose to implement this approach:

- a) *The organization is deploying analytics and data warehousing for the first time.* In this situation, the analytic platform may be the initial step in building out a traditional EDW environment. One of the 8 customers interviewed for the study fit into this category. The customer chose an analytic platform that enabled the organization to start with a small data warehouse appliance, but grow, via a set of scalable offerings, to provide a large EDW system that can support multiple business areas.
- b) *The organization does not have sufficient budget, time, or resources to incorporate the data into an existing EDW.* In this situation, an analytic platform offers the promise of deploying this so-called *independent data mart* solution at a lower cost and a shorter time to value. In the future, depending on business need, the data in the data mart may be integrated into an EDW. Many companies have learned from experience, however, that independent data marts may save time and money in the short term, but may prove more costly in the long term because data marts have a tendency to proliferate, which creates data consistency and data integration issues. As a result, many experts have a negative view of the independent data mart approach.
- c) *The organization needs to support extreme processing where it is unnecessary or impractical to incorporate the data into an EDW.* Six of the customers interviewed for this research report match this scenario. Depending on business need, the independent analytic solution may acquire data from an EDW to augment the analyses and may also replicate the processing results back into an EDW. Some independent analytic solutions may be experimental in nature or may only exist for a short period of time to fulfill certain short-term analytic needs.

The first 2 reasons, or scenarios, just outlined are well understood because they are normally a part of the traditional data warehousing life cycle. The extreme processing scenario, however, is relatively new and represents the biggest potential for business growth and exploitation of analytics. It is important, therefore, to look at extreme processing in more detail.

There are several factors driving the need for extreme processing. The first is the growth in data volumes, number of data sources, and types of data. As we noted earlier, many organizations are now generating tens of terabytes of data per day. For these organizations, it is becoming impractical, or even impossible, for cost, performance, or data latency reasons to load certain types of data (high volume web event data, for example) into an EDW. In some cases it may not even be necessary. The application may involve data that only has a useful lifespan of a few days or weeks. Note, however, that these latter types of applications do not preclude the analytic results, or scored or aggregated data from being stored in an EDW for use by other analytic applications.

Another factor driving extreme processing is the nature of the analytical processing itself. BI users are becoming more knowledgeable and more sophisticated in their use of analytics. They want to analyze detailed data as well as aggregated data. They are also building more complex analyses and more advanced predictive models. There is also an increasing demand by these users for enabling ad hoc analyses, in addition to the more traditional predefined reports and analyses provided by IT.

Extreme data coupled with extreme analytical processing leads to the need for high performance and elastic scalability. In data-driven companies, many analytic applications are mission critical, and reliability and high availability are therefore also of great importance. Given constantly changing business requirements and data volumes, the analytic platform in these situations needs to support flexible hardware growth and also be easy to build, manage, expand, and if necessary, tear down and replace. These extreme needs require a new approach to data warehousing, and, in our opinion, this is the sweet spot for new and evolving analytic platforms. These analytic solutions do not replace the traditional data warehousing approach—they extend it by enabling extreme processing.

To use the term *independent data mart* to describe the underlying data store and data management system supporting extreme analytic application processing misrepresents this new breed of applications and the business benefits it can provide. Perhaps a more suitable term would be an *extreme analytic platform*.

### Use Case 3: Filtering, Staging, and Transformation of Data

The objective of this use case is to exploit the parallel processing power of the analytic platform's ADBMS to perform data filtering and transformation. This approach is particularly useful in environments involving high volumes of data and/or a wide variety of data sources and types of data. Note that the NoSQL software (Hadoop with MapReduce, for example) discussed earlier is a strong competitor to this approach.

The processing of the data in this use case is typically done using an ETL approach where the:

- *Extract* step collects and filters source data
- *First load* step loads the filtered data into a set of temporary staging tables in the ADBMS
- *Transform* step does the required transformation and integration of the filtered data
- *Second load* step loads the transformed data into the ADBMS or a remote DBMS for analytic processing

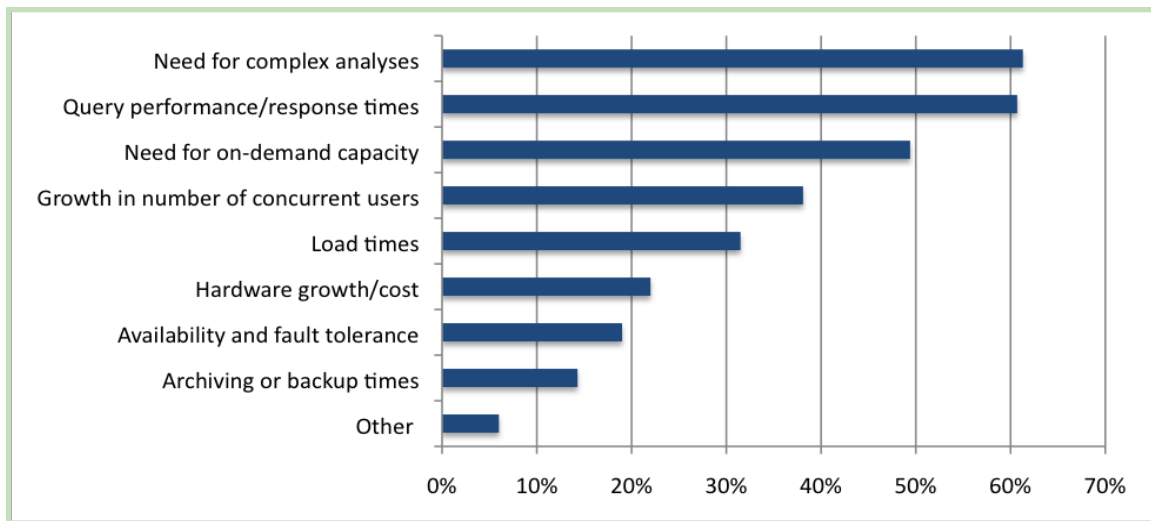
Some 42% of survey respondents stated they were using an analytic platform for the filtering, staging, and transformation of data. One of the customers interviewed for this report was using an ETL approach with an extreme analytic platform. The business users in this case were able to analyze both the detailed data and the aggregated results from the ETL processing.

The analytic processing performed in this use case supports data transformation and aggregation, rather than the creation of business analytics. One use of this scenario is to transform less well-structured data into a more usable format. Textual data (web pages, blog pages, unstructured log files, for example) is a strong candidate for this type of transformation.

This use case also offers an alternative to using extreme processing. Instead of loading high-volume detailed data into an extreme analytic platform, an intermediate system is used to filter and/or aggregate the detailed data so that it is practical and cost-effective to load it into an EDW. Of course, the downside of this approach is that information is lost in the filtering and aggregation processing.

We can see from the research study survey results and customer interviews that analytic platforms are being used to support all three of these use cases. In our survey, we also asked organizations what

circumstances caused them to move to, or consider, an analytic platform for supporting these use cases. The results are shown in Figure 7.



**Figure 7: What Issues Led You to Use an Analytic Platform?**

The top 5 reasons were: need for complex analyses (61%), query performance (61%), on-demand capacity (49%), growth in user audience (38%), and load times (32%). These results clearly demonstrate that cost is not the main driving force behind using an analytic platform. This was confirmed by our customer interviews. Only one company indicated that reducing costs was the key reason for replacing its existing EDW solution with a new analytic platform.

The top five reasons given in the survey for using an analytic platform, however, do have an indirect relationship to cost. Many of the customers interviewed for this report were building extreme analytic solutions, and the reasons they gave for choosing any given analytic platform all matched one or more of the top five results from the survey. Most of these companies were deploying applications that couldn't be built before. This was either because the application couldn't provide the required performance no matter how much hardware was employed or because the amount of hardware required to achieve acceptable performance was cost prohibitive.

Cost and performance are therefore related, but the key takeaway from the results is that analytic platforms provide cost-effective solutions that extend, rather than replace, the existing data warehousing environment. They enable applications that simply could not be built before.

## Getting Started

Success in business is an elusive thing: solving today's question opens new possibilities and new questions for tomorrow. Few categories of technology remain as consistently at the top of CIO planning—"What have you done for me lately?" is the typical question managers are asked. Success for the procurement staff is a signed contract; for business analysts the challenge is greater. The following are some thoughts on how to ensure that the platform selected works today and will continue to grow and evolve as your needs do.

There are also best practices for the implementation of your analytics strategy that leverage the tools you acquire and following them will make success more likely. The vendors and users we interviewed offered some valuable insights and we include them here, together with a quick review of the success factors that make the difference.

### Selecting the Right Platform

For analytic platform selection, there is only one place to start: understanding the analytics planned. The use cases in this study hint at the possibilities, but they also make it clear that there is enormous variety: in the skills and preferred tools of the users, the business problems being tackled, the types of analysis required, the latency of the data, and the users' volumes. Will standard reporting be enough? Is ad hoc analysis with drill-down and slice-and-dice operations required? Does internal and external data need to be combined? Will the analyses involve data mining and predictive model building? Will temporary analytic data stores be set up, processed, and then torn down frequently? What are the current and future data volumes and number of users? Know these answers before you begin. Otherwise, making vendor choices is a hit-and-miss process that is likely to lead to project failures.

A vendor POC is the single most vital part of the product selection process. No selection process, however complex, can substitute for the one critical element: testing on your data, with your queries, on the hardware and software platform you plan to use, with the number of concurrent users that matches the expected usage patterns. As you decide who should be on your short list for actual tests, here are some key aspects to consider as you draw up your requirements.

- **Getting the data in and keeping it available.** Ensure that you can load data at the speed you need to absorb it from the sources you expect to use and that any filtering, transformation, and distribution/partitioning you will need to do is supported. Assess the tools offered for design—how complex are they? Do they optimize for your expected queries automatically? Are changes possible without taking the system down for long periods? How does the system provide backup and recovery? How does it assure availability if failures occur?
- **Working with your languages and tools.** Unless you want to train your users and developers extensively, look for products that work with the tools you are familiar with. Consider users beyond the usual suspects internally—you may hope to involve more departments, add more skills, and answer more questions than you have before.
- **Supporting your toughest questions.** If you did your homework, you should know the tough questions that need to be answered and the features that are required to answer them. Complex joins, multipass processing, sophisticated statistics and mashups can make or break products—from the most mature on down.

Other aspects have to do with deployment specifics—don't ignore basics like the interface to your storage hardware, the speed of the interconnects, the level of pre-integration provided across the software stack. Be sure that setting up test or development systems is no more complex than you are comfortable with; analytics is an increasingly iterative process. Explore the possibility of doing such work in the cloud—can the vendor support that? How difficult would it be to move from completed testing in the cloud to production on your hardware?

Finally, the POC trial process is a great indication of the vendor's ability to support you. If the personnel involved don't seem knowledgeable, problems take a long time to resolve, and/or setup seems to take

a long time, you must consider what it will be like when the check has been signed and the purchase made. Assess what services are available to you for design, training, and support. And be sure to leave some surprises. **Do not** conduct your trial on prearranged queries and analyses only. Stress test workloads that mirror your expected ones in number of users, volumes of data, and other processes running if there will be any.

## Conclusions

Analytic platform adoption is strong, and accelerating. Hundreds of millions of dollars being spent with new vendors represent the most significant spending shift in database systems in a decade. Customers are satisfied; the promises made are being kept. As prospective purchasers test difficult problems in POC exercises, shortlisted products are proving equal to the tasks—beating incumbent classic RDBMSs. The right selection process involves understanding the likely analytical workloads, data volume and types, and numbers of concurrent users—all should be tested. POCs separate winners from losers: often, some candidates can't get it done at all.

Support for open source data integration and BI tools, columnar data storage, and a complete hardware and software-packaged solution are not yet top of mind for purchasers. Conversation with early adopters and survey data show that query performance, support for complex analytics, and on-demand capacity are. As analytic platforms become mainstream, however, it's likely that ease of installation and support and aggressive data compression strategies will begin to grow in importance.

In 2010, analytic platform offerings from DBMS leaders Oracle, Microsoft, and IBM entered the market, and this development should drive increased awareness and growth. The analytic platform will drive billions of dollars in revenue in the next decade, and transform expectations about the ability to use data to improve business results.

## Appendix: Detailed Survey Results

**Q1:** *We define an analytic platform as: "An integrated and complete solution for managing data and generating business analytics from that data, which offers price/performance and time to value superior to non-specialized offerings. This solution may be delivered as an appliance (software-only, packaged hardware and software, virtual image) and/or in a cloud-based software-as-a-service form." Do you agree with this definition?*

Value	Count	Percent %
Yes	209	93.7%
No	14	6.3%
<b>Total Responses</b>		<b>223</b>

**Q2:** *Are you using or planning to use an analytic platform?*

Value	Count	Percent %
Already using	99	44.4%
No plans	55	24.7%
Planning/creating specifications	42	18.8%
Evaluating products for a defined project	27	12.1%
<b>Total Responses</b>		<b>223</b>

**Q3:** *How much historical data do you keep online (in non-archival form) for analysis?*

Value	Count	Percent %
More than 3 business years or 12 quarters	95	42.6%
3 business years (or the past 12 complete quarters) or less	55	24.7%
1 business year (or the past 4 complete quarters) or less	51	22.9%
90 days (or the past business quarter) or less	22	9.9%
<b>Total Responses</b>		<b>223</b>

**Q4:** *Do you routinely perform business analysis on data that is not maintained in an RDBMS?*

Value	Count	Percent %
No	105	47.1%
Yes, with hand-coded programs	86	38.6%
Yes, with packaged tools	32	14.3%
<b>Total Responses</b>		<b>223</b>

**Q5:** Which of the following use cases [architectural models] are being deployed for your analytic platform? Check all that apply.

Value	Count	Percent %
Enterprise data warehouse (EDW)	115	68.5%
Data staging area for a data warehouse	70	41.7%
Dependent data mart	70	41.7%
Independent data mart or data store	65	38.7%
Other	7	4.2%
<b>Total Responses</b>	<b>168</b>	

**Q6:** Which of the following issues led you to add an analytic platform? Check all that apply.

Value	Count	Percent %
Need for complex analyses	103	61.3%
Query performance/response times	102	60.7%
Need for on-demand capacity	83	49.4%
Growth in number of concurrent users	64	38.1%
Load times	53	31.5%
Hardware growth/cost	37	22%
Availability and fault tolerance	32	19%
Archiving or backup times	24	14.3%
Other	10	6%
<b>Total Responses</b>	<b>168</b>	

**Q7:** How much raw data are you managing on your analytic platform? (Raw data is the source data loaded into a data store before adding indexes, aggregate tables, materialized views and/or cubes built from the raw data.)

Value	Count	Percent %
1 to 10 terabytes	67	39.9%
Less than 1 terabyte	62	36.9%
11 to 20 terabytes	18	10.7%
21 to 100 terabytes	12	7.1%
Greater than 100 terabytes	9	5.4%
<b>Total Responses</b>	<b>168</b>	

**Q8:** *How much data are you managing on your analytic platform after loading, tuning, enhancing, and compressing the raw data?*

Value	Count	Percent %
Less than 1 terabyte	73	43.5%
1 to 10 terabytes	68	40.5%
21 to 100 terabytes	12	7.1%
11 to 20 terabytes	10	6%
Greater than 100 terabytes	5	3%
<b>Total Responses</b>		<b>168</b>

**Q9:** *How many concurrent users do you need your analytic platform to support?*

Value	Count	Percent %
Less than 20	61	36.3%
21 to 100	53	31.5%
101 to 1,000	38	22.6%
Greater than 1,000	16	9.5%
<b>Total Responses</b>		<b>168</b>

**Q10:** *What data sources are used to feed your analytic platform? Select all that apply.*

Value	Count	Percent %
Structured RDBMS data	132	78.6%
Structured file data	91	54.2%
Structured legacy DBMS data	78	46.4%
XML data	75	44.6%
Packaged application data	49	29.2%
Unstructured file data	40	23.8%
Weblogs	39	23.2%
Event or message data	37	22%
Data from enterprise service bus or web service	29	17.3%
Rich media data	11	6.5%
Other	11	6.5%
<b>Total Responses</b>		<b>168</b>

**Q11:** Please rate the following features that were/are important in acquiring or planning your analytic platform environment?

	Very Important	Somewhat Important	Not Very Important
Query performance	77.4%	19.6%	3.0%
Reduced load times	58.3%	32.7%	8.9%
Good administration tools	54.2%	39.9%	6.0%
Fault tolerance and high availability	54.2%	36.9%	8.9%
Integration into IT environment	53.6%	38.7%	7.7%
Easy scaling & hardware upgrades	45.2%	41.7%	13.1%
Support for commercial DI/BI tools	43.5%	43.5%	13.1%
In-database processing	33.9%	47.6%	18.5%
Workload management	31.0%	56.0%	13.1%
Data compression	31.0%	50.0%	19.0%
Support for open source DI/BI tools	27.4%	49.4%	23.2%
Support for cloud computing	26.2%	33.9%	39.9%
Columnar storage	25.6%	56.0%	18.5%
In-memory data	24.4%	49.4%	26.2%
Support for big data (>100TB)	23.8%	40.5%	35.7%
Complete H/W & S/W package	18.5%	47.6%	33.9%
MapReduce support	10.7%	44.6%	44.6%

**Q12:** Has your analytic platform project met your expectations?

Value	Count	Percent %
Partially	116	69%
Fully	36	21.4%
No	16	9.5%
<b>Total Responses</b>		<b>168</b>

**Q13:** *What industry is your company in?*

Value	Count	Percent %
Computer Services/Consulting	35	15.7%
Financial/Banking/Insurance/Real Estate/Legal	32	14.3%
Computer software/hardware/technology manufacturer	22	9.9%
Business Services/Consulting	17	7.6%
Government	16	7.2%
Communications/Telecom Supplier	15	6.7%
Education	15	6.7%
Health/Health Services	13	5.8%
Retail/Wholesale	10	4.5%
Manufacturing/Industry (non-computer related)	9	4%
Other (please specify)	9	4%
Service Provider (ASP, ESP, Web hosting)	4	1.8%
Manufacturing consumer goods	4	1.8%
Travel/Hospitality/Recreation/Entertainment	3	1.3%
Aerospace	3	1.3%
Other	25	11.2%
<b>Total Responses</b>	<b>223</b>	

**Q14:** *How many employees (worldwide) are in your company?*

Value	Count	Percent %
1 to 49	46	20.6%
1,000 to 4,999	35	15.7%
100,000	25	11.2%
5,000 to 9,999	23	10.3%
100 to 249	18	8.1%
10,000 to 24,999	17	7.6%
50,000 to 99,999	13	5.8%
500 to 999	13	5.8%
250 to 499	12	5.4%
25,000 to 49,999	11	4.9%
50 to 99	10	4.5%
<b>Total Responses</b>	<b>223</b>	

**Q15:** *On whose behalf are you completing the survey?*

Value	Count	Percent %
Complete company	73	32.7%
Business department	50	22.4%
Consulting client	39	17.5%
Business division	31	13.9%
Other	30	13.5%
<b>Total Responses</b>		<b>223</b>

**Q16:** *Please tell us where you and your company are located.*

Value	North America	Europe	Asia/Pacific	Latin America
Where are you located?	59.6%	13.5%	18.4%	8.5%
Where is your corporate HQ?	65.9%	15.2%	13.0%	5.8%

## Aster Data Overview and Business Description

Aster Data ([www.asterdata.com](http://www.asterdata.com)) provides an MPP database management system with an integrated analytics engine to enable cost-effective management of large data volumes and rich analytics on large data sets. Its unique coupling of SQL with the MapReduce analytics framework enables high performance parallel processing. Much of its R&D has been dedicated to the enablement of high performance, scalable queries and advanced in-database analytic processing in its flagship “Data Analytics Server” called Aster Data *nCluster*. The company believes that processing full data sets together with application logic on one platform is a key requirement for analytic data warehouses, enabling deeper insights from the data, more precise models, and higher performance for more real-time, mission-critical applications. It sees demand for a single platform for different types of applications; its customers want to put all their data in one place. Seventy percent (70%) of the data it sees is not living in the EDW; Aster Data says its customers are looking for a place to aggregate it. In-database processing eliminates the latency inherent in shipping data across the network to a processing tier. Analysts and developers can access MapReduce through standard SQL because of Aster’s SQL-MapReduce framework.

Privately held, Aster Data shipped its first product in 2007 and in the three years since has amassed several dozen customers, mostly in high-end deals targeted at large volumes of data like those at comScore, Barnes & Noble, Akamai, and MySpace. Its key markets include digital media, financial services, government, and retail. Most of its business is in the U.S. and Europe, and some has come through a strong partnership with Dell, including joint marketing, sales, and an appliance-based offering.

### Architecture

Aster Data *nCluster* is an MPP database with a hybrid row- and column-oriented data storage architecture and an integrated analytics engine that runs on commodity hardware. It targets large—terabyte and petabyte scale—databases that ingest many rich data types. *nCluster* uses four node types: Queen nodes for coordination, Worker nodes for distributed analytical processing, Loader nodes for high throughput data loading, and Backup nodes for massively parallel backup. The MPP architecture makes linear scaling possible. Node specialization for loading is also a critical attribute; Aster Data claims 8 terabyte per hour load rates and also supports trickle loading. Data compression rates are competitive; comScore (customer case study follows) reports achieving 8:1 compression for one of its systems.

Its universal computation layer can source data from row- or column-oriented tables and perform both SQL-based and MapReduce-based analysis. Aster Data *nCluster* has a unique multitiered architecture which enables task isolation and the ability to scale each tier incrementally as needed to meet workload requirements for query processing, data loading, and backup.

*nCluster* provides fault tolerance with replication, automatic failover, failure heuristics, and clustered backup to prevent unplanned downtime due to hardware or software failures. Its dynamic workload manager allocates processing and compute resources to in-progress transactions, allowing administrators to change priorities in real time with rule-based prioritization (preadmission control) and dynamic resource allocation and reallocation. The Aster Data Management Console offers GUI-based control of physical architecture, process-level inspection, and resource control.

On premises (appliance or otherwise) is not the only deployment option; Aster Data works with Amazon, AppNexus, and Terremark to offer cloud deployment. Sharmila Mulligan, EVP of worldwide marketing, notes, “Some verticals have strong preferences—in the federal and financial services space, the preference is on premises—often with specific hardware and even specific operating system releases. Web companies, other than the largest, want to go to the cloud for minimal implementation and maintenance costs.”

### **Analytic Functionality**

The MPP environment enhances performance dramatically by running analytic functions on each node in the cluster, parallelized across nodes and even across cores. Aster Data’s customers have both application and ad hoc needs that require powerful analytics—which can be either difficult to express with only SQL or perform poorly (due to complex joins.). The patent pending SQL-MapReduce framework enables powerful procedural programmability integrated tightly with standard ANSI SQL for business analyst simplicity and tight BI ecosystem fit. Aster Data’s focus on ease of development for rich analytic applications is very visible in its Eclipse-based integrated development environment (IDE), Aster Data Developer Express. It automates the creation of packages into which analysts import their Java code, enables local testing on the client, and provides one-click push down of the analytic application to the *n*Cluster servers. The IDE is freely available for download at Aster Data’s website, [www.asterdata.com](http://www.asterdata.com).

The SQL-MapReduce analytics framework is designed to offer parallelized execution while preserving isolation from database operations to ensure high availability and high performance. Unlike stored procedures and UDFs, Aster Data’s in-database processing is polymorphic, meaning that functions are not tightly bound to only one table or data structure in the database—they are reusable. Users of the functions can be prompted for values at runtime without requiring developers to rewrite the code. For complex analytics, this can greatly reduce time to value as models, segmentations, or other analytic approaches are tested. *n*Cluster scales MapReduce analysis with the data; Aster Data’s dynamic workload manager can automatically redistribute and reallocate workloads. Developer Express also speeds analytic application development by automatically generating MapReduce code so the developer does not need to learn MapReduce but rather can leverage existing Java and SQL skills.

### **Differentiation**

Aster Data *n*Cluster is differentiated by both the MPP database platform and the mechanisms for enabling advanced analytics. *n*Cluster’s independently tiered architecture ensures that reads, writes, backup and data loading always occur in parallel. Each tier in the architecture can be scaled independently and runs on standard, commodity hardware, so that *n*Cluster can scale with both data and budget. One-click administration is provided for scaling and ensures that *n*Cluster automatically installs and configures new commodity hardware nodes without requiring system downtime. This means “always-on” availability.

System availability is further ensured with background replication processes that store replicated data across servers so that in the event of failures (which will inevitably happen from time to time with commodity hardware), the system automatically activates the replicas, allowing queries to complete as if nothing happened. Dynamic workload management ensures highly predictable performance and guaranteed service levels for complex mixed workloads, including reads, writes, and loads. Fine-grained policy controls allow administrators to reallocate CPU and storage resources based on in-progress transactions to meet usage requirements.

The layer “above” the engine and the interfaces is where Aster Data focuses its differentiation story—on analytics. To speed development of custom and advanced analytic functions, the Developer Express IDE lets developers write in their language of choice, from SQL to Java, C, C++, C#, Perl, Python, .NET, and R. Aster Data also provides Analytic Foundation, a suite of ready-to-use SQL-MapReduce functions intended to accelerate development. The functions are delivered in two ways. “Business analyst-ready functions” provide deep algorithms (linear regression or K-means) and enable them through SQL so that the user can simply set parameters and turn the function loose on the data, across all nodes in the cluster in parallel. The result can operate like a prompted query—add a user interface and prompt the user. “Power-user functions” are more programmer oriented, highly specific building blocks to be put into programs created by coders.

Analytic Foundation, bundled with the enterprise edition of the *nCluster*, offers functions including statistical analysis, clustering (including K-means), time series, sessionization, graph and market basket analysis, data transformation, Monte Carlo simulation, histograms, linear algebra, geospatial and text processing (including “unpack,” which takes weblogs apart for use in analytics).

### **Partnerships**

Aster Data, a young company, has already built strong partnerships with technology vendors such as Dell and HP; analytics application providers such as SAS; BI tool vendors SAP BusinessObjects, IBM Cognos, and MicroStrategy; and data integration players like Informatica and Pentaho. Its reseller relationships with Carahsoft in the federal market and with Amazon for prospects interested in cloud deployments have also been key.

On the analytics front Aster Data partners with SAS, Fuzzy Logix, Cobi Systems for geospatial development; Cloudera and Impetus around integrating Hadoop; and Ermas for in-database SAS and R.

### **How Should Customers Start, and What Matters Most?**

Aster Data believes customers need to focus on several key areas: performance, scalability, richness, and ease of development. Concentrating on one at the expense of the others will limit flexibility going forward. It’s important to envision what types of analytics are needed—advanced reporting, operational apps, ad hoc analysis, or highly interactive analysis. It’s preferable to point to a wide range rather than just optimize for one type.

### **Future/Road Map Exploitation of Trends**

Aster Data believes that a platform for diverse data types consolidated into a single platform will be in high demand—combining relational, non-relational, and unstructured data with diverse analytic applications that access data directly from multiple storage engines via a SQL and NoSQL interface, including SQL-MapReduce as a universal analytics framework for all data types. It continues to focus on delivering more out-of-the-box, prepackaged analytics modules (e.g., time-series analysis, cluster analysis, graph, statistical analysis, etc.) to ease development of rich analytic applications. To further ease development of rich analytic applications on Aster Data’s platform, Developer Express will encompass higher level functionality to move beyond automated glue code creation and code completion to areas such as drag-and-drop code component.

## Aster Data Customer Case Study: comScore

### Company Background

comScore ([www.comscore.com](http://www.comscore.com)) is a leading provider of digital marketing intelligence. With approximately 2 million worldwide panelists under continuous measurement, the comScore panel utilizes a sophisticated methodology that is designed to accurately measure people and their behavior in the digital environment. This information enables comScore's clients to better understand, leverage, and profit from the rapidly evolving world of the web-based and mobile computing. Its services are used by more than 1,600 organizations in over 40 countries.

comScore is a public traded company (NASDAQ: SCOR) founded in August 1999. Among the areas of digital marketing that comScore was first to measure are e-commerce (2001), search (2002), ad networks (2005), global Internet audience measurement (2005), online video (2005), widgets (2007), and mobile Internet browsing (2007). It has some 900 employees in U.S., Europe and Asia. The company is headquartered in Reston, Virginia.

For this case study, we interviewed Will Duckworth, Vice President of Software Engineering.

### The Business Problem

In late 2008, comScore wanted to extend its services beyond those enabled by using data gathered from its 2 million panelists. It wanted to introduce a new census-based service that uses visitor traffic data collected from the web sites of its partners. For these services, comScore needed to be able to load in excess of 18 billion new rows of data a day for analysis. It also needed a configuration that could scale easily as it added new partners to the service. comScore decided, therefore, to look for a new analytic solution for managing the census-based data.

### The Analytic Platform Solution

Part of the solution chosen was the Aster Data *nCluster* analytic platform using in-house developed data integration software for the hourly update of the analytic store, which consumes about 1.1 terabyte of compressed data per day. In-depth data analysis is done using complex SQL queries developed by business analysts; query results are often exported to flat files and Microsoft Excel. The hardware environment consists of a 70-node Dell server configuration running CentOS Linux with a storage capacity of 7.2 terabytes of data per node. comScore is expanding this configuration to support higher than anticipated data volumes and the ability to store 90-days of history data.

Aster Data *nCluster* was selected after running a proof-of-concept trials with several vendors to evaluate performance and functionality. In addition to meeting comScore's performance requirements, other reasons for selecting Aster Data *nCluster* were support for standard SQL syntax, the ability to use MapReduce functions and embed them into the DBMS, and the capability to easily and cost-effectively expand the server and disk configuration to meet data growth. Will Duckworth also noted that another reason for choosing Aster Data was that, "The company was easy to interact with at all levels of management and had a similar culture to that of comScore."

The initial deployment of the application will support business analysts, but is expected to grow to include development and quality assurance analysts at comScore as it goes into full production. These analysts are primarily interested in exploring detailed data produced during the last 24 hours to 30 days, but data needs to be maintained in the system for 90 days to monitor trends and to enable the next generation of product development.

Transformed and aggregated output data from MapReduce processing in comScore's Hadoop environment is also brought into the Aster Data *nCluster* system. This is done using file transfers at present, but this approach will be replaced in the future by Aster's Hadoop Connector feature.

### Implementation Considerations

The initial 10-node system was installed in early 2009 and was operational less than 48 hours after it was delivered. "It's easy and fast to add additional hardware as we need it," said Duckworth.

Hardware delivery and software certification delays created some initial problems. This shortage was due to using hardware that was not in general release and not certified by Aster Data. "These delays sometimes caused us to suddenly make big changes to the disk configuration, which can lead to a packet storm of data moving around the system as it automatically rebalances itself," said Duckworth. "We are loading data 24 hours a day, and we need to load an hour's worth of data in an hour in order to stay on schedule. It's much better to gradually upgrade the system to avoid potential disturbances," he added.

Perhaps one of the most significant challenges was determining what data business analysts required in the analytic store for analysis. "We needed to have good communications with business users about their needs," said Duckworth. "Given the amount of data involved, changes in data structure, such as adding a new data column, can have a major impact on the system in terms of partitioning, indexing, etc. It also affects historical information and historical comparisons. These types of changes have to be carefully planned," Duckworth added. "The Aster platform enables a range of new types of powerful queries to be created. These again need to be carefully thought out and planned for in order to achieve the full benefits that Aster offers."

comScore has upgraded to the latest release of Aster Data *nCluster* in part because of its new workload management feature. Duckworth said he was impressed with this feature and the other new capabilities in Version 4.5. comScore particularly wants to leverage the power of the enhanced MapReduce functionality.

### Benefits

comScore is currently building out its census-based solution and the Aster Data platform. "We are pleased with the lower cost to scale out the configuration and are planning to take full advantage of the benefits that the Aster MapReduce capabilities offer," said Duckworth.

### Success Factors

The two main success factors in the comScore project were the performance of the Aster Data software and the support the organization received from the vendor. Aster Data gave the comScore implementers good access to all levels of management, which enabled them to resolve issues quickly and have improvements made to the software to meet their requirements.

### Summary

The analytic platform use case for comScore was an analytic application and associated data store for a new line of business. Key selection criteria were lower hardware costs, fast load times, the ability to scale to support increasing data volumes, and good analytic processing performance. Other factors that played an important role in choosing Aster Data *nCluster* were support for complex SQL queries and the ability to program and embed application processing into the DBMS using MapReduce. A good working relationship with the software vendor was also important.

## About the Authors

**Merv Adrian**, Principal at IT Market Strategy, has spent 3 decades in the information technology industry. As Senior Vice President at Forrester Research, he was responsible for all of Forrester's technology research for several years, before returning to his roots as an analyst covering the software industry and launching Forrester's well-regarded practice in Analyst Relations. Prior to his Forrester role, Merv was Vice President and Research Manager with responsibility for the West Coast staff at Giga Information Group. Merv focused on facilitating collaborative research among analysts, and served as executive editor of the monthly Research Digest and weekly GigaFlash. He chaired the GigaWorld conference (and later Forrester IT Forum) for several years, and led the jam band, a popular part of those events, as a guitarist and singer.

**Colin White** is the president of DataBase Associates Inc. and founder of BI Research. As an analyst, educator and writer he is well known for his in-depth knowledge of data management, information integration, and business intelligence technologies and how they can be used for building the smart and agile business. With many years of IT experience, he has consulted for dozens of companies throughout the world and is a frequent speaker at leading IT events. Colin has written numerous articles and papers on deploying new and evolving information technologies for business benefit and is a regular contributor to several leading print- and web-based industry journals. For ten years he was the conference chair of the DCI and Shared Insights Portals, Content Management, and Collaboration conference. He was also the conference director of the DB/EXPO trade show and conference.