

Специализированные решения для DW и бизнес-анализа

Обзор специализированных решений для хранилищ данных (Data Warehouse – DW) и приложений бизнес-анализа, представленных на региональном рынке.



Константин Клячкин — заместитель директора департамента системных решений, компания IBS.



Максим Исаев — менеджер по развитию бизнеса департамента системных решений, компания IBS.

Введение

Современный успешный бизнес — это бизнес, который чутко и быстро реагирует на запросы и изменения рынка, может не только предлагать индивидуальные услуги клиентам, но и проактивно изменять их с учетом потребностей клиентов в режиме онлайн. В этой связи в последнее время резко повысился интерес к специализированным DW-решениям (СПХД — специализированные платформы для хранилищ данных и бизнес-анализа), который обусловлен, прежде всего, двумя факторами. Это:

- возросший динамизм бизнеса с требованиями его адаптивности к быстро меняющимся условиям рынка в высококонкурентной среде;

- снижение порога ценовой доступности подобных решений. Последнему обстоятельству способствует также и развитие облачных сервисов, среди которых предлагаются и BI-сервисы.

Можно выделить 3 основные этапа развития СПХД:

- **этап 1:** до 2005 г. — рынок представлен всего несколькими ключевыми игроками — NCR/Teradata, IBM и др. Стоимость СПХД достаточно высока;
- **этап 2:** с 2005 г. по 2010 г. — на рынке появляется множество игроков — Teradata, IBM, Oracle, HP (Neoview Platform), DATAlegro, Netezza, Sybase, Kognitio и др. Стоимость СПХД становится намного демократичней;
- **этап 3:** с 2010 г. — на рынке появляются СПХД для всех секторов бизнеса — от малого до большого. Сам рынок снова укрупняется. Специализированные DW-решения становятся одним из ключевых направлений для всех крупных мировых вендоров — Teradata, IBM, Oracle, EMC, HP, Microsoft, SAP/Sybase и др. Появляются специализированные решения для бизнес-анализа не только для структурированных, но и для неструктурированных данных. Объемы DW могут достигать десятков петабайт. На рынке появляются облачные BI-сервисы, не требующие капитальных вложений (в ближайшие 1–2 года подобные серви-

сы должны появиться и в России, прим. ред.).

На рис. 1. представлены основные типы пользователей прикладных BI-систем с точки зрения их функций в компаниях и количества. Это позволяет понять задачи, решаемые на каждом уровне и дает возможность оценить требования, предъявляемые пользователями на каждом уровне к ХД.

Первое поколение хранилищ данных (ХД) решало в основном стратегические задачи — было ориентировано в основном на верхне-уровневый анализ деятельности компании на основе регламентных отчетов и KPI. Количество пользователей было невелико, как и объемы анализируемых данных. Второе поколение решало тактические задачи — например, задачи маркетинга (повышение спроса), логистики (сокращение затрат), финансов (сокращение издержек) и т.п. Это требовало более частого обновления данных и более глубокой работы с ними. Помимо руководства компаний, системами начали пользоваться линейные менеджеры и аналитики. Российские компании лишь относительно недавно начали переходить на этот уровень, в то время как во всем мире это — стандарт довольно давно.

Следующее поколение хранилищ должно быть ориентировано, прежде всего, на динамизм современного бизнеса. Условно их можно назвать — операционные ХД. Стра-

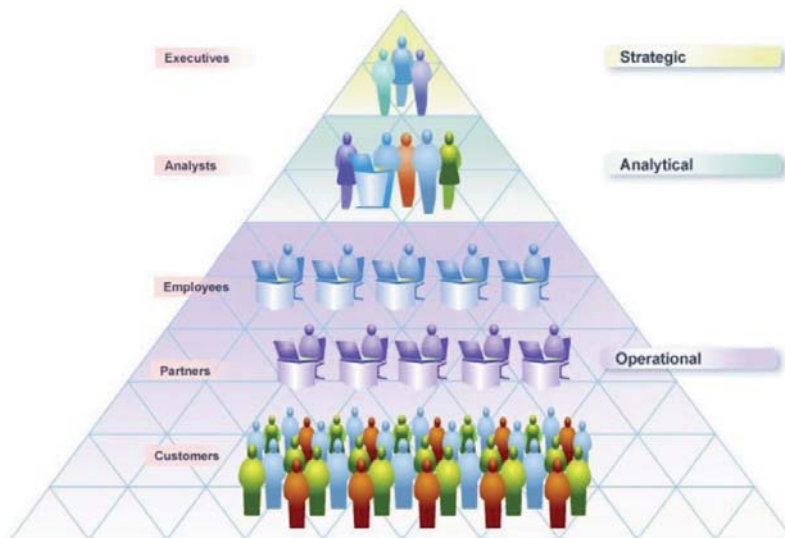


Рис. 1. Основные типы пользователей прикладных BI-систем с точки зрения их функций в компаниях и количества.

тегия и тактика многих компаний меняется настолько стремительно, в зависимости от ситуации на рынке, что основным требованием к ХД становится работа в режиме реального времени. В то же время, фокус анализа сместился на продукты, процессы, клиентов и т.п. В результате, на основании подобных анализов компании должны получить возможность управлять взаимоотношениями с любыми клиентами в любом месте и в любое время, прогнозируя их поведение на шаг вперед. Но решение подобных задач не отменяет задач стратегических и тактических, поэтому новое поколение ХД должны были обеспечивать решение всей совокупности задач. Среди российских компаний едва ли единицы имеют такие системы.

Реализация всех подходов шла с учетом технологических возможностей каждого периода времени. Сначала это были реплики оперативных БД, поэтому и были придуманы ХД с витринами и кубами, для того чтобы обходить ограничения традиционного оборудования и ПО. Уже в то время начались попытки интеграции ПО и оборудования. В общем, они оказались довольно успешными. Что касается традиционных подходов, то объемы данных и количество пользователей росло быстрее и в результате столкнулись с ограничениями.

С развитием e-бизнеса, глобализацией экономики все большее количество компаний будут испытывать потребности в системах, обеспечивающих, помимо высокой скорости доступа к данным, такую же скорость принятия сложных управленческих решений на основе всей совокупности накопленной информации, хранимой в базе данных, превышающей многие терабайты, и уже петабайты в условиях необходимости одновременной обработки запросов от сотен и тысяч клиентов в секунду. В качестве примеров таких отраслей можно привести следующие:

- **финансовые учреждения** должны правильно и, по-возможности, в режиме реального времени оценивать все существующие риски, например, при кредитовании клиентов или проведении торговых сделок на бирже, а операторы call-center должны в момент звонка клиента понимать какой дополнительный продукт или сервис ему можно было бы еще предложить;
- **компании розничной торговли** должны иметь единое интегрированное представление каждого заказчика во всем множестве каналов продаж, его “историю” — и все это в разрезе продуктовых линеек, чтобы обеспечить оптимальное обслуживание с учетом персонализированных предложений;
- **телекоммуникационные компании:** бизнес становится все более конвергентным. Нужно понимать поведение клиентов, быстро реагировать на малейшие его изменения, снижать расценки для уходящих клиентов и предлагать дополнительные услуги/сервисы, персонализированные для каждого клиента с учетом его интересов. Кроме того, законодательства обязывают их хранить огромные объемы исторических данных.

Специализированные DW-решения имеют многолетнюю историю, начало которой было положено компанией Teradata в конце

90-х. Отличительными особенностями специализированных DW-решений являются:

- наличие параллельной архитектуры, поддерживаемой как аппаратным обеспечением, так и специализированным ПО для распараллеливания обработки запросов, что, например, позволяет обрабатывать тысячи запросов в секунду при глубине анализа до нескольких десятков и сотен последних транзакций по каждому запросу;
- высокая масштабируемость с возможностью наращивать производительность практически линейно — до сотен узлов, а объем хранилища данных — до многих петабайт;
- автоматизация выполнения многих рутинных операций без необходимости вмешательства какого-либо обслуживающего персонала. Например, возможность одновременного выполнения сервисных процедур по актуализации данных (индексация, агрегирование, подсчет суммарных значений отдельных колонок таблиц, актуализация данных СПХД с онлайн-новыми данными продуктивных систем и др.) вместе с выполнением запросов пользователей снимает и большую проблему по тюнингу этих процедур, упрощает управление БД и резко повышает доступность СПХД (приближая ее к уровню 24x7);
- наличие встроенной оптимизации запросов. Спецификой современных СПХД является то, что они должны не только поддерживать работу с очень большими базами данных, но и обработать чрезвычайно сложные запросы к ним, которые просматривают, соединяют и одновременно выполняют вычисления с большим количеством таблиц в системе. При этом число таблиц, используемых в запросе, может достигать до 30 и более. Для примера: при наличии в запросе 15 отношений между таблицами возможно более чем 1,3 трлн различных объединений. Наличие встроенной в ядро СПХД функциональности по обработке таких запросов позволяет уменьшить число таких отношений, тем самым значительно снижая время их выполнения;
- возможность поддержания сложной смешанной BI-нагрузки без снижения производительности. За счет встроенного архитектурного параллелизма СПХД можно одновременно выполнять сложную загрузку данных и их синхронизацию, одновременно удовлетворяя требованиям SLA с тысячами пользователей и

обеспечивая мониторинг и управление всех параллельных задач.

Рассмотрим, как решают вопросы реализации СПХД ряд ключевых представителей этого рынка.

Teradata

Основанная в 1979 г. в Кремниевой Долине Teradata является родоначальником хранилищ данных и построила первый в мире программно-аппаратный комплекс ХД — Teradata Database Computer (DBC/1012), который уже в 1986 г. был признан Fortune Magazine “Продуктом года”. С самого начала при разработке своих решений Teradata была ориентирована на схемы массовой параллельной обработки данных — Massive Parallel Processing (MPP). Три принципа являются основополагающими при построении СПХД Teradata.

1. Параллельная обработка данных за счет равномерного распределения данных по всему дисковому пространству.

Основными компонентами, обеспечивающими параллелизм, являются Parsing Engine (PE) и Access Module Processor (AMP). PE обрабатывает SQL запросы, обеспечивая наилучший план их выполнения, и посылает команды по обработке данных на AMP в соответствии с полученным планом. AMP отвечает за управление только определенной частью дискового пространства. Таким образом, один AMP “видит” только некоторую порцию каждой таблицы, хранимой в БД. Поэтому в случае запроса на массовую обработку, за счет равномерного распределения данных каждой таблицы, нагрузка распределяется между всеми AMP в системе, обеспечивая максимально возможную производительность.

2. Инновационные решения по обработке данных, среди которых следующие два.

Teradata Virtual Storage (начиная с версии СУБД Teradata 13.0) представляет дисковое пространство как единый пул по хранению данных, контролирует выделение дискового пространства для конкретного AMP, что упрощает управление дисками.

Управление нагрузками (workload management). Система позволяет гибко распределять и приоритезировать нагрузку процессов использования и обработки данных в зависимости от времени суток, количества сессий, сложности запроса, загруженности ресурсов (CPU/memory).

3. Безграничная линейная масштабируемость. BYNET версии 4 способна объединить работу до 4096 узлов.

Табл. 1. Семейство целевых платформ Teradata для построения хранилищ данных.

	Data Mart Appliance 5XX	Extreme Data Appliance 1XXX	Data Warehouse Appliance 2XXX	Extreme Performance Appliance 4XXX	Active Enterprise Data Warehouse 6XXX
Цель	Тестирование, разработка, небольшие витрины данных	Аналитика на огромных объемах данных, для новых типов данных	ХД или витрины данных уровня департамента	Экстремальная производительность для операционного анализа	ХД уровня предприятия, активное ХД уровня предприятия
Масштабируемость	до 12 ТВ	до 200 PB	до 343 ТВ	до 18 ТВ	до 114 PB
Сегмент применения	Аналитика уровня отдела, сервер начального уровня	Аналитический архив, Deep Dive Analytics	Strategic Intelligence, поддержка принятия решений, быстрое сканирование	Operational Intelligence, низкий объем, высокая производительность	Активная нагрузка, обновления в реальном времени, тактические и стратегические запросы

Перечисленные основные принципы доступны во всех решениях Teradata. В зависимости от целей и задач, поставленных бизнесом, Teradata предлагает 5 классов систем для построения ХД предприятия (табл. 1).

Teradata Data Mart Appliance — это начальное, полностью интегрированное решение для витрин данных и ХД уровня департамента, объемом до 12 Тбайт, построено на базе одного узла SMP и СУБД Teradata, может использоваться в том числе для целей тестирования и разработки.

Teradata Data Warehouse Appliance 2XXX (2650/2690) — платформа для построения хранилищ и витрин данных в отдельных подразделениях крупных корпораций до 343 Тбайт. Решения данного семейства оптимизированы для обеспечения максимальной скорости доступа к данным и используются в основном для поддержки принятия стратегических решений. В октябре 2011 г. Teradata анонсировала 5-е поколение этих решений — 2690, которое имеет удвоенную производительность и утроенную емкость по сравнению со своим предшественником.

Teradata Active Enterprise Data Warehouse 6XXX (6650/6680) — флагманский продукт корпоративного класса. Может включать в свою конфигурацию от 1 до 4096 узлов и хранить до 114 Пбайт пользовательских данных. Этот вид ХД способен обеспечить решение одновременно как тактических, так и стратегических задач бизнес-аналитики, используя гибкое управление нагрузками, допуская одновременное выполнение массовых загрузок/выгрузок данных, генерацию аналитической отчетности, выполнение ad-hoc запросов для решения конкретных задач, потоковую обработку данных. Активное хранилище данных подразумевает необходимость решения тактических задач, где необходимо обновление данных в режиме, близком к реальному, высокая производительность и доступность системы.

В марте 2011 г. Teradata приобрела компанию Aster Data Systems, также поставляющую решения ХД, основным преимуществом которых является инновационный подход к анализу больших объемов данных (Big Data). Анализ данных выполняется не на транзакционном, а на событийном уровне, который является намного более детальным. Данная модель позволяет анализировать такие данные, как клики пользователя на сайте при веб-серфинге и при наборе корзины покупок в интернет-магазине, оценить наиболее часто используемые функции мобильного устройства (телефона/планшета), поведение клиентов в социальных сетях, анализировать машинные логи и логи приложений (например, количество резких торможений при управлении автомобилем, что позволяет страховым компаниям понять стиль вождения клиента и предложить скидку на страховку).

Из достижений Teradata на российском рынке можно отметить бурный рост числа ее клиентов за последние 3 года — до 12 — один из самых высоких в регионе для данного рынка. Среди них: Сбербанк, ВТБ-24, МТС, Ростелеком и др.

Oracle

Oracle Exadata Database Machine

Корпорация Oracle — один из наиболее активных игроков на рынке СПХД. Свой выход на него Oracle ознаменовала в конце сентября 2008 г. решением Oracle Database Machine / Exadata — специализированной платформой для хранилищ данных на базе Oracle 11.1, предназначенной как для хранения данных, так и частичной или полной обработки SQL-запросов на уровне СХД. Exadata представляет собой кластер параллельно работающих узлов на базе стандартных серверов. Параллелизм обработки SQL-запросов к такому ВІ-хранилищу данных достигается за счет ряда технологий: “умного” ПО — Exadata Storage Server Software, отвечающего за правильное хранение данных (данные распределены по всем дискам) и способного задействовать в нужный момент вычислительные ресурсы аппаратного кластера серверов хранения. Таким образом, каждый сервер выполняет свою часть работы, а результаты в дальнейшем консолидируются сервером БД.

Первая версия Exadata была ориентирована, прежде всего, на работу с ВІ-приложениями, которым приходится обрабатывать таблицы БД размером от сотен мегабайт до нескольких терабайт, где часто необходимо выполнять полное сканирование таблиц. В качестве классических примеров можно назвать ВІ-системы, отчетные системы и им подобные. Уже в первой своей реализации Exadata позволила получить ускорение обработки ВІ-запросов до 50–70 раз (в среднем до 16–28 раз) при проведении тестирования для компаний розничной торговли и телеком-компаний (*по данным Oracle, прим. ред.*). Однако результаты могут быть и более скромными для ряда задач и/или неоптимальной настройке.

В настоящее время на территории России имеется несколько сертифицированных центров, проводящих тестирование Exadata непосредственно на задачах пользователя.

Exadata позволяет балансировать и приоритизировать различные нагрузки как между различными группами/классами пользователей/приложений внутри одной базы, так и между базами данных, гарантируя при этом заданный (в соответствии с SLA) уровень выделения ресурсов ввода/вывода.

Вторая версия Oracle Exadata (*продажи с сентября 2009 г., прим. ред.*) уже в полной мере поддерживала и бизнес-критичные OLTP-приложения, как по скорости доступа, так и по поддержанию доступности данных.

В июне 2011 г. Oracle объявила о сертификации приложений SAP для использования с Exadata. Теперь приложения SAP, основанные на архитектуре SAP NetWeaver 7.x, такие как SAP ERP и SAP Business Warehouse (BW), которые получили сертификаты на совместимость с базой данных Oracle Database 11g Release 2, могут разворачиваться с использованием Oracle Exadata Database Machine X2-2 и X2-8.

Третье расширение возможностей Exadata произошло в октябре 2011 г., когда Oracle объявила о полной оптимизации бизнес-приложений Oracle для машины баз данных Oracle Exadata Database Machine и машины связующего программного обеспечения Oracle Exalogic Elastic Cloud. Благодаря интеграции удалось существенно повысить производительность таких приложений, как: Oracle E-Business Suite, Oracle PeopleSoft, Oracle JD Edwards EnterpriseOne, Oracle Siebel CRM, Oracle ATG Commerce Suite и Oracle Supply Chain Management, но и дополнить их встроенной ВІ-функциональностью. В частности, интеграция с Exadata и Exalogic отразилась на производительности приложений следующим образом:

- *Oracle E-Business Suite*: операции при управлении кадрами и закупками осуществляются в 8, транзакции “от заказа до оплаты” — в 3, а бухгалтерские операции — в 7 раз быстрее;
- *Oracle PeopleSoft*: интеграция повысила скорость и масштабируемость приложений PeopleSoft для управления финансами и персоналом. В финансовых приложениях PeopleSoft журнальные записи стали вноситься в бухгалтерские книги в 5 раз быстрее, а расчет заработной платы ускорился на 40%. Масштабируемость приложений PeopleSoft для управления персоналом стала в 10 раз выше, а время отклика для функций самообслуживания сократилось в 3 раза;
- *Oracle JD Edwards EnterpriseOne*: интеграция ускоряет пакетную обработку на 33% и обеспечивает самое низкое время отклика в сравнении с результатами всех когда-либо проводившихся тестов производительности полнофункциональных систем для управления ресурсами предприятия;
- *приложения семейства Oracle Supply Chain Management*: интеграция значительно сокращает время на планирование цепочки поставок, а анализ спроса в POS-системах теперь занимает 45 минут вместо шести часов.

Oracle Exalytics Business Intelligence Machine и Oracle Big Data Appliance

Следующий шаг в развитии направления СПХД был сделан Oracle в октябре 2011 г. с объявлением Oracle Exalytics Business Intelligence Machine и Oracle Big Data Appliance (*начало продаж — 1 кв. 2012 г., прим. ред.*).

Oracle Exalytics Business Intelligence Machine — первый в отрасли оптимизированный программно-аппаратный комплекс, который позволяет компаниям расширить использование ВІ-средств, переходя от отчетов и информационных панелей к моделированию, планированию, составлению прогнозов и предиктивному анализу. Приложения для планирования могут теперь работать в масштабе всего предприятия с более короткими и более точными циклами планирования. Exalytics позволяет достигать очень высокой производительности (*по заявлениям разработчика, прим. ред.*) при аналитической обработке данных. С по-

явлением Exalytics стал доступен высокоскоростной анализ данных для многих тысяч пользователей мобильных устройств, т.е. любому человеку независимо от его местонахождения.

В качестве технологических платформ в Oracle Exalytics используются:

- серверы Oracle Sun Fire с ОЗУ объемом 1 Тбайт и процессорами Intel Xeon E7-4800 с общим числом ядер 40;
- Oracle BI Foundation Suite, включающий Oracle Business Intelligence Enterprise Edition и Oracle Essbase со средствами оптимизации производительности и средой расширенной визуализации для интерактивного анализа любых данных и поддержки любых классов пользователей;
- решение Oracle TimesTen In-Memory Database for Exalytics, разработанное на основе поставляемой Oracle и проверенной на практике высокопроизводительной реляционной системы управления базами данных, обрабатываемой в оперативной памяти, расширенной и оптимизированной для бизнес-анализа.

Существующие системы отчетов и информационных панелей, разработанные на основе Oracle BI Enterprise Edition и Oracle BI Applications, могут выполняться на Oracle Exalytics без каких-либо изменений.

Oracle Exalytics – это открытое решение для гетерогенных ИТ-сред, способное осуществлять доступ и анализировать данные, размещенные в любых реляционных базах данных, как Oracle, так и других поставщиков, многомерных OLAP-серверах и в неструктурированных источниках данных. Среди поддерживаемых платформ – IBM DB2, Microsoft SQL Server, Netezza, SAP Business Information Warehouse и Teradata.

Для организаций, которым требуется еще более высокая производительность, Oracle Exalytics предлагает коннектор InfiniBand, предназначенный для работы с машиной баз данных Oracle Exadata Database Machine.

Внутренние тесты, сравнивающие Oracle Exalytics с типовым программным обеспечением для бизнес-анализа и готовыми аппаратными решениями, продемонстрировали прирост производительности до 20 раз при составлении реляционных OLAP (ROLAP) отчетов и при использовании информационных панелей, и до 79 раз – при многомерном OLAP (MOLAP) моделировании. Резкий рост производительности многомерного моделирования достигался при размещении кубов Oracle Essbase в оперативной памяти. Все горизонтальные и отраслевые приложения Oracle BI, основанные на платформе Oracle BI Foundation Suite, также показывают повышенную производительность.

Oracle Big Data Appliance

Oracle Big Data Appliance (BDA) дает возможность эффективно использовать т.н. “большие данные”, которые генерируются блогами, социальными сетями, интеллектуальными счетчиками/датчиками и

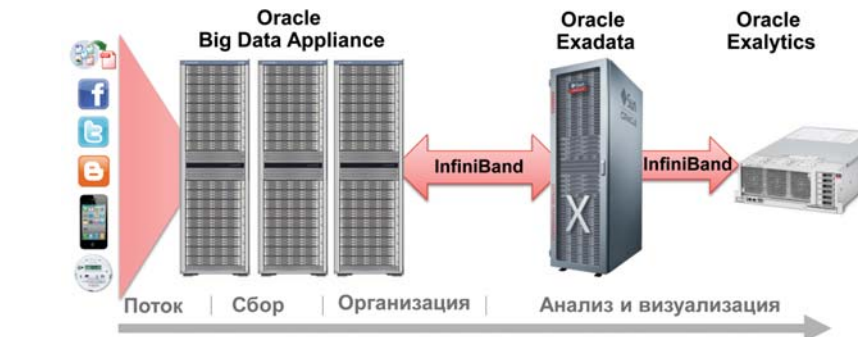


Рис. 2. Oracle Big Data Appliance вместе с Oracle Exadata Database Machine и Oracle Exalytics Business Intelligence Machine представляет собой полное решение для аналитической высокоскоростной обработки всех типов данных.

веб-сайтами, телеком-оборудование и т.п., которые сегодня, зачастую, не попадают в ХД и не доступны для обработки и анализа.

Oracle Big Data Appliance вместе с Oracle Exadata Database Machine и Oracle Exalytics Business Intelligence Machine представляет собой полное решение (рис. 2) для высокоскоростного анализа всех типов данных: 1) данных OLTP-приложений; 2) “больших данных”, доступных для онлайн-анализа после пакетной обработки средствами BDA и размещения полученной структурированной информации в хранилище данных (Exadata); 3) онлайн-данных.

Тема “больших данных” активно стала обсуждаться сравнительно недавно, когда стало понятно, что появился большой слой данных, который требовал собственных технологий хранения и обработки/анализа, “не вписывающихся” в предложения, представленные на рынке.

Наибольший вклад в бурный рост “больших данных” внес интернет и такие проекты, как: LinkedIn, Facebook, Digg, Google+, Amazon, Ebay, Yahoo и др. В качестве особенностей “больших данных” можно назвать следующие:

- пониженные требования к целостности данных, например, потеря нескольких старых логов пользователей при посещении ими интернет-магазина не приведет к каким-либо последствиям;
- необходимость хранить огромное количество данных, измеряемое петабайтами;
- ценность данных, в отличие, например, от OLTP-данных, падает очень медленно. Поэтому необходимость прямого доступа к ним может измеряться годами и большими периодами. Например, анализ изменения спроса на отдельные группы товаров в ритейловых сетях может проводиться в течение многих лет. Более того, могут представлять интерес какие-то новые запросы по новым признакам, изначально не заложенные в систему. В этой связи крайне ценным становится хранение всей изначально собираемой информации, а не ее “выжимки”;

- данные плохо структурированы или вообще не структурированы, более того, структура данных может меняться.

В результате для реализации подобных систем стали предлагаться кластеры на базе стандартных дешевых серверов, допускающих масштабирование до сотен и тысяч узлов. Так, в 2006 г. появился проект Hadoop с открытым кодом (проект Apache Software Foundation) для хранения (на базе распределенной файловой системы HDFS с возможностью автоматической обработки отказа узлов и перераспределения данных) и пакетной обработки в массивно-параллельном режиме “больших данных”.

В качестве узлов для Oracle BDA используются серверы Sun X4270 M2 (48 Гбайт памяти, два шестиядерных процессора Intel Xeon (L5640), 12 HDD – 2 Тбайт 7200 RPM SAS, два порта 40 Гбит/с InfiniBand). Таким образом, BDA, состоящий из четырех стоек, содержит 60 узлов (18 узлов в стойке, всего 864 ядер) и 1,7 Пбайт для хранения данных. При этом система имеет неограниченную масштабируемость.

В состав ПО BDA входят: Oracle Linux 5.6, Java Hotspot VM, Apache Hadoop Dis-

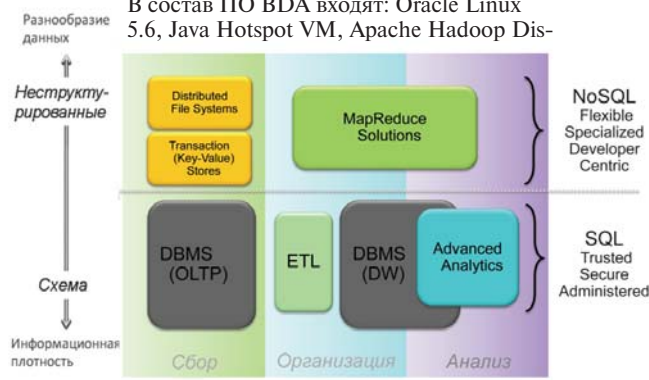


Рис. 3. Разделение технологий при сборе, структуризации и анализе данных.

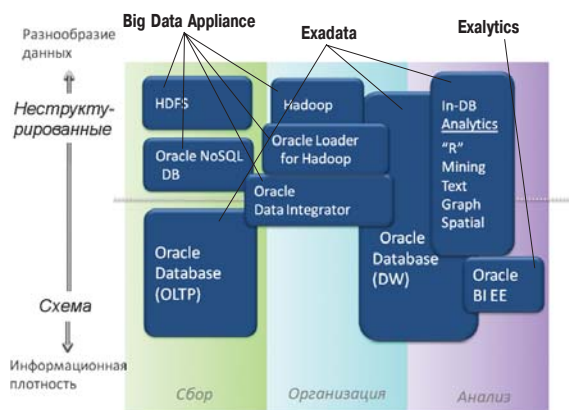


Рис. 4. Разделение технологий между специализированными BI-системами Oracle.

tribution v0.20.x, R Distribution, Oracle NoSQL Database Enterprise Edition, Oracle Data Integrator Application Adapter for Hadoop, Oracle Loader for Hadoop, оптимизированное ПО (Hadoop HDFS, HBase and Hive, NoSQL Database и реплики).

Разделение технологий при сборе, структуризации и анализе неструктурированных и структурированных данных, а также разделение технологий между специализированными BI-системами Oracle представлены на рис. 3,4.

IBM

В настоящее время в портфеле IBM — четыре семейства решений класса СПХД. Два из них — IBM Netezza и IBM Smart Analytics System — представлены на рис. 5. Эти решения — программно-аппаратные комплексы. В отличие от них, IBM InfoSphere Warehouse — это чисто программное решение для создания хранилищ данных собственными силами.

IBM Netezza — масштабируемые системы с массовым параллелизмом, позволяющие обрабатывать петабайтные объемы данных при множестве одновременно выдаваемых запросов высокой сложности. Компанию Netezza IBM приобрела еще в ноябре 2010 г., однако продвигать и продавать BI-решения в России на основе ее технологий стала только в октябре 2011 г. По заявлениям производителя, отличительной особенностью IBM Netezza является то, что это решение практически не требует никакой настройки самого хранилища и после загрузки — оно полностью готово к работе.

Линейка IBM Netezza состоит из трех продуктов (рис. 6). На высокопроизводительную аналитику ориентировано решение IBM Netezza 1000 (табл. 2) с масштабируемостью до 10 стоек (емкость несжатых данных — 320 Тбайт, емкость с компрессией — 1,28 Пбайт, число S-Blades — 120).

IBM Netezza 1000 — это специализированное решение для управления хранилищами данных на основе стандартных компонент, обеспечивающее интеграцию баз данных, серверов, систем хранения и функций расширенной аналитики в единую систему, удобную для управления. Оно предназначено для быстрого и углубленного анализа данных, объем которых измеряется десятками терабайт или петабайтами. IBM Netezza 1000 содержит множество серверов Snippet Blades, или S-Blades, на которых выполняются фрагменты кодов запросов SQL (так называемые "сниппеты") и сложные аналитические процессы. Серверы S-Blades — это интеллектуальные узлы обработки, которые образуют систему с массовым параллелизмом. Каждый S-Blade представляет собой независимый сервер, имеющий в своем составе многоядерные процессоры Intel, уникальные многофункциональные FPGA IBM Netezza и гигабайты ОЗУ — сбалансированные и работающие параллельно, обеспечивающие высокую производительность.

Каждому жесткому диску в составе IBM Netezza соответствует выделенный процессор, таким образом, при загрузке данных происходит равномерное распределение нагрузки по всем дискам и процессорам. При добавлении стоек коэффициент распараллеливания увеличивается пропорционально их числу. Своей высокой производительностью семейство IBM Netezza 1000 обязано архитектуре асимметричной обработки данных с массовым параллелизмом (Asymmetric Massively Parallel Processing, AMP™), в которой blade-серверы IBM и дисковые накопители интегрированы с патентованными инструментами фильтрации данных IBM на основе программируемых логических матриц типа FPGA. Такое сочетание обеспечивает очень высокое быстродействие при выполнении запросов в условиях сверхсложных разнотипных рабочих нагрузок для поддержки десятков тысяч пользователей бизнес-аналитики и хранилищ данных (по заявлениям разработчика, прим. ред.).

Для средств интеграции и загрузки данных, а также инструментов анализа и визуализации, СПХД IBM Netezza 1000 выглядит как обычная реляционная СУБД, доступ к которому осуществляется посредством стандартных интерфейсов ODBC, JDBC и OLE DB. Перекладывание аналитических задач на IBM Netezza не должно вызывать затруднений, так как в составе идет встроенная платформа аналитики IBM Netezza Analytics.

Встроенное ПО — IBM Netezza Analytics поддерживает множество аналитических инст-

Табл. 2. Состав IBM Netezza 1000 в разных конфигурациях.

Параметры	Системы в одной стойке			Системы в нескольких стойках	
	1000-3	1000-6	1000-12	2 стойки	3+ стойки
IBM Netezza 1000	1000-3	1000-6	1000-12	2 стойки	3+ стойки
Стойки	1	1	1	2	3-10
S-Blades	3	6	12	24	число стоек x 12
Число ядер процессора	24	48	96	192	число стоек x 96
Объем пользовательских данных, Тбайт (без сжатия)	8	16	32	64	число стоек x 32

рументов и языков программирования; поставляется она с библиотекой аналитических функций для баз данных, которые выполняют аналитические операции параллельно, скрывая сложность параллельного программирования от разработчиков.

Недавно IBM анонсировала поддержку Netezza для мэйнфреймов в качестве back-end устройств для решения задач бизнес-анализа.

Вторым решением в области СПХД предлагаемым компанией IBM является — IBM Smart Analytics System (ребрендинг решения IBM InfoSphere Balanced Warehouse, которое IBM ввела с марта 2008 г., прим. ред.) построенный на оптимизированных компонентах, в состав которых входят серверы, СХД, элементы коммутации, специализированное ПО. Масштабирование осуществляется блоками. Параллелизм обработки запросов поддерживается как внутри блоков, так и при добавлении новых блоков.

IBM Smart Analytics System имеет два принципиальных отличия от Netezza:

- возможность поддержки смешанных нагрузок — OLTP-запросов и аналитических запросов;
- для IBM Smart Analytics System требуется настройка для каждого типа нагрузки. Т.е. достижение сопоставимой производительности с Netezza возможно, но для этого потребуются усилия квалифицированного персонала и время. Если нагрузка не детерминирована и меняется, то накладные затраты и время на настройку могут быть значительны. Хранилище данных в IBM Smart Analytics System организовано по типу традиционной реляционной СУБД, т.е. необходимо поддерживать индексацию, агрегирование и др.

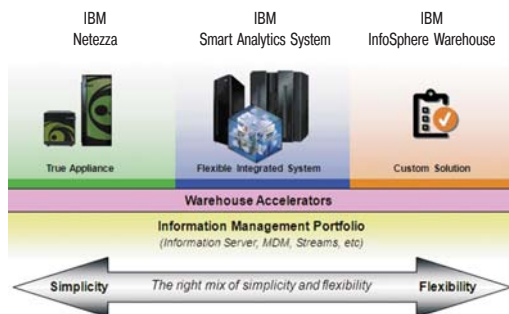


Рис. 5. Семейство специализированных решений IBM для хранилищ данных BI.

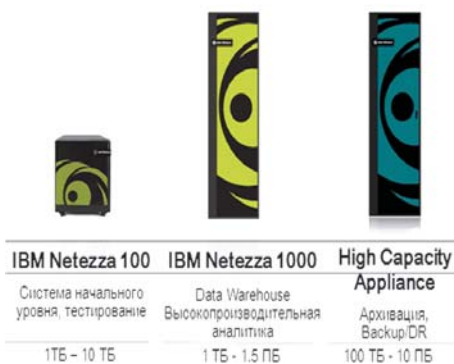


Рис. 6. Линейка продуктов IBM Netezza.

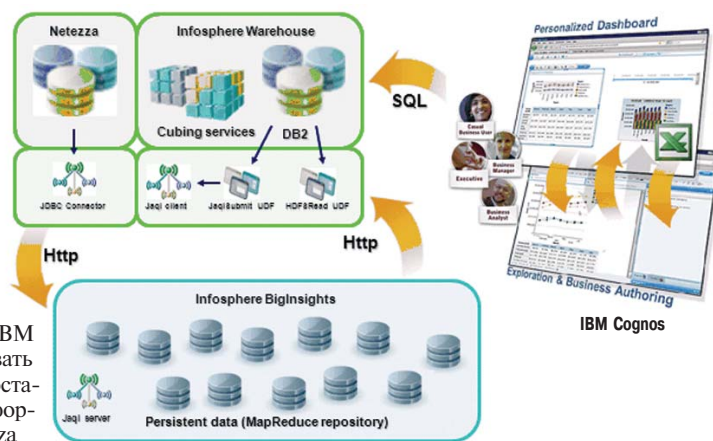


Рис. 7. Интеграция решений IBM Netezza, IBM Smart Analytics System, IBM BigInsights и IBM Cognos позволяют осуществлять глубокий аналитический анализ как структурированных, так и неструктурированных данных, дополняя друг друга.

Третье семейство решений, которое IBM предлагает для продвижения направления СПХД — IBM BigInsights. Оно лишь косвенно относится к данному тренду, поскольку предназначено для упрощения развертывания и управления Hadoop-кластерами и ориентировано на работу с т.н. “большими данными”.

Все три отмеченных решения могут быть интегрированы, работать в составе одного ЦОД и дополняться системами визуализации и управления данными (рис. 7).

Последнее — четвертое семейство решений IBM для СПХД — специализированные решения для SAP HANA, обеспечивающие ускорение работы BI-приложений за счет максимальной обработки “in-memory”. Основное ограничение на применение таких систем — узкий класс задач, хранилища данных, которых могут размещаться в памяти сервера.

Заключение

Итак, чем же СПХД могут быть интересны для бизнеса и за счет чего они эффективнее традиционных решений?

Первое преимущество — в первую очередь, это повышение производительности аналитических систем при масштабировании DW до петабайтного уровня, что позволяет выполнять предикативные запросы с широкими возможностями по моделированию развития ситуации в будущем с максимальным приближением к реальному времени, а также поддерживать сложную аналитику в онлайн — в прямом и переносном смысле “не отходя от кассы”. Это дает возможность развивать совершенно новые направления в бизнесе и по-новому подходить к его организации вследствие появившихся возможностей.

За счет чего же происходит оптимизация и повышение скорости работы аналитических систем? Прежде всего, благодаря тому, что Hardware Appliances протестированы и сконфигурированы оптимальным образом самими производителями ПО для BI (например, Oracle, IBM, Microsoft, SAP/Sybase и др.). Поскольку они разрабатываются для решения определенных задач, то у производителей есть возможность за счет устранения ненужного функционала направить высвобождающиеся вычислительные мощности на повышение быстродействия при решении основных задач.

Кроме того, за счет использования специализированных технологий (например, Netezza, Teradata) удается поддерживать высокий уровень производительности аналитических систем для разных BI-нагрузок без каких-либо перенастроек, что становится особенно важно в быстроменяющейся среде.

Ряд СПХД обеспечивают интеграцию продуктивных БД и хранилищ данных практически в реальном времени, что еще в большей степени способствует принятию правильных решений в то время, как только возможные признаки кризисных явлений появляются.

Что касается второго преимущества — снижения совокупной стоимости владения, то оно является следствием возможности экономить и/или увеличить прибыль за счет:

- поддержки и ПО, и оборудования одним вендором, что позволяет существенно сократить сроки решения проблемы;
- высокой производительности, благодаря которой появляется возможность работы нескольких ресурсоемких приложений на одном сервере (или нескольких, если используется кластер). В ре-

зультате сокращения времени простоя и динамической балансировки нагрузки в пиковые часы обеспечивается более рациональное использование ресурсов;

- поставки решения полностью готового к работе (нет от всех вендоров, прим. ред.) без необходимости дополнительных настроек, что сокращает затраты на управление и администрирование, а также уменьшает сроки запуска систем в промышленную эксплуатацию;
- более эффективного механизма хранения и работы с данными, благодаря чему снижаются потребности в дисковых массивах;
- уменьшения затрат на энергоресурсы;
- возможности решения новых задач, принятия более быстрых и правильных решений. Это, в свою очередь, позволяет компании наращивать конкурентоспособность и увеличивать прибыль.

И, конечно, при оценке различных решений надо помнить: лучшим из них является не самое громоздкое и дорогое, а грамотно спроектированное.

Компания IBS как лидер российского рынка консалтинга и информационных технологий имеет наработанный проектный опыт и команду высококвалифицированных специалистов по большинству из этих обсужденных продуктов. Специалисты компании IBS обладают всей необходимой компетенцией и опытом для предоставления широкого спектра услуг заказчикам в случае принятия ими решений о начале использования специализированных программно-аппаратных комплексов для решения своих задач.

**Константин Кляцкин,
Максим Исаев,
компания IBS**

NetApp® Open Solution for Hadoop

Ноябрь 2011 г. — Компания NetApp (NASDAQ: NTAP) объявила о выпуске преконфигурированного готового к внедрению решения NetApp® Open Solution для платформы Hadoop, которое позволит клиентам получить максимальную отдачу от внедрения Hadoop-кластера благодаря высокой гибкости и производительности и добиться снижения ТСО.

На сегодняшний день производится больше данных, характеризующихся все большим объемом и сложностью, чем когда-либо ранее — от данных с терминалов розничной торговли и информации с мобильных телефонов до научных и маркетинговых исследований. У корпоративных клиентов теперь есть возможность использовать новые способы управления и анализа этих данных для стимулирования инноваций, принятия более взвешенных решений и достижения успешных результатов, которые помогут создать конкурентные преимущества в бизнесе.

NetApp и компания Cloudera, Inc. объединили свои усилия для распространения дистрибутива Cloudera Distribution,

включающего платформу Apache Hadoop (CDH) и Cloudera Enterprise, абонентскую услугу, в которую входит Cloudera Support и управляющее ПО для Hadoop, вместе с решением NetApp Open Solution for Hadoop, что позволит ускорить внедрение и использование Apache Hadoop заказчиками.

NetApp Open Solution for Hadoop представляет собой комплексный, готовый к внедрению модульный кластер Hadoop для корпоративного внедрения платформ Hadoop. Данное решение имеет следующие особенности и преимущества:

- быстрое внедрение и простое масштабирование;
- высокая общая производительность кластера;
- снижение затрат за счет предоставления технологии самовосстановления, RAID и глобальных резервных дисков;
- снижение общей стоимости покупки и эксплуатации.

Решение NetApp Open Solution for Hadoop предоставляет клиентам доступ и инструменты для анализа имеющихся данных и извлечения прибыли из своих информационных активов. Клиенты могут использовать уже имеющиеся у них серверы и базы данных для анализа очень больших объемов данных и получения

с течением времени максимальной окупаемости инвестиций в свои данные.

“NetApp давно поддерживает открытые стандарты и предоставляет своим клиентам необходимые им гибкость и эффективность для управления результатами своих данных. Наша стратегия для Hadoop будет точно такой же, мы будем предоставлять клиентам выбор для внедрения Hadoop — от системы корпоративного класса до решений Apache Hadoop с открытым кодом», — отметил Рич Клифтон, страший вице-президент и генеральный директор подразделения внедрения технологий и решений NetApp. — Cloudera усиливает наше корпоративное решение NetApp Open Solution for Hadoop, вместе мы упрощаем корпоративное внедрение Hadoop, облегчаем управление и сокращаем стоимость внедрения”.

CDH является основным инструментом для внедрения на предприятиях платформы Apache Hadoop и ее производственного использования. Компании могут управлять полным жизненным циклом своих систем Apache Hadoop, поскольку CDH обеспечивает полную видимость кластеров Hadoop. Он также позволяет автоматизировать текущие системные изменения, необходимые для поддержания и улучшения качества операций.