

## Уровни технологий для обработки Big Data

**Н**а сегодняшний день рынок технологий для обработки Больших Данных предлагает самые различные решения по своему составу, которые можно разделить на пять категорий. Самая базовая категория, первая ступень, – это аппаратное обеспечение под Big Data. Следующая категория – файловые системы, которые строятся поверх аппаратного обеспечения, а также файловые системы особого класса, которые хорошо масштабируются и подходят под решения Big Data. Третий уровень – это платформы аналитики для Больших Данных. Четвертый уровень – уже готовые базы данных. И наконец, пятый уровень – готовые решения, которые включают в себя аналитическую платформу, базу данных и аппаратное обеспечение для начала работы с Big Data. Расскажем про каждый уровень по порядку.

Начнем с базового класса – **аппаратного обеспечения**. С точки зрения архитектуры хранения в нем все выстраивается гораздо проще, чем мы привыкли видеть в больших корпоративных решениях. Дело в том, что традиционное хранение данных и традиционные решения для хранения данных, используемые в enterprise-классе (обычно это консолидированные системы с общим доступом – NAS-системы), не обеспечивают должной масштабируемости и производительности для задач Big Data. Единственный способ получить требуемые характеристики – это создать распределенную платформу, в которой будет множество вычислительных узлов и множество узлов хранения. Обычно для этого используется архитектура Shared Nothing (без общих элементов). В этом случае за счет массивного параллелизма, то есть распределения задач по множеству узлов, удается получить требуемую производительность и масштабирование.

Эти результаты достигаются за счет идеологии, которая заключается в том, что в связи с таким большим объемом данных существует необходимость “приблизить” процессор к данным, а не данные к процессору, как это происходило ранее. Для этого пришлось фактически вернуться к очень простым решениям, придумав алгоритм, который способен распределять нагрузку по большому количеству процессоров и контролировать результаты на выходе. Этот алгоритм, созданный компанией Google, называется MapReduce. Одной из его имплементаций является всем известный Hadoop. Благодаря этим разработкам на рынке были созданы

условия, при которых для построения аналитических платформ для обработки Больших Данных используются наиболее простые (с точки зрения архитектуры и технологий) решения по хранению. Обычно это самые простые серверы с внутренними дисками либо DAS-хранилищем, то есть то, чем пользовалась индустрия хранения еще в 60-80-х годах прошлого века. Поэтому готовые платформы или платформы с открытой архитектурой, которые поставляются разными производителями для обработки Больших Данных, являются довольно простыми устройствами и вообще похожи друг на друга. Компания Dell производит различное аппаратное обеспечение для облачных вычислений. Это x86-серверы и недорогие, но быстрые хранилища класса DAS, которые используются для подключения к вычислительным узлам. В принципе такие архитектуры очень похожи на то, что можно увидеть в нижеупомянутых комплексных предложениях от различных производителей.

Вторая категория предложений – это **файловые системы**. Весь описанный выше объем для хранения данных, собранный из открытого аппаратного обеспечения на x86-платформе с дисками с прямым подключением, для начала работы нужно собрать в логическое хранилище. Для этого используются либо файловые системы в составе пакета Hadoop (HDFS), либо готовые промышленные файловые системы с проприетарными файловыми системами, которые имеют интеграцию с Hadoop. Таким образом, у пользователя есть выбор – собрать хранилище самостоятельно, используя аппаратное обеспечение открытой архитектуры и бесплатные Open Source-пакеты, либо приобрести готовое решение, которое предоставляет распределенную, масштабируемую, высокопроизводительную файловую систему. Здесь только важно понимать, что это не будет готовым “решением Big Data”, это лишь хранилище – решение одной из множества задач, которые необходимо реализовать в таком проекте.

Следующий, третий, уровень – это, соответственно, покупка или самостоятельное развертывание **аналитической платформы**. Здесь можно использовать как уже вышеупомянутый пакет Hadoop, в том числе его коммерческие дистрибутивы, или же его альтернативы (например, разработку Greenplum компании EMC). Также стоит отметить, что Hadoop не является целостным программным обеспечением – на самом деле это целый пакет программных разработок. Соответствен-

**Dell Cloudera Hadoop Solution**

**Benefits**

- #1 distribution of Apache Hadoop
- Cloudera Enterprise
  - Cloudera Manager
  - Professional Services
  - Support & Upgrades
  - Training & Certification

**Components**

- Making Open Source Hadoop Enterprise Ready
- CDH3 Enterprise
- Dell-developed Crowbar Software
- Dell PowerEdge C, 12g Servers
- Dell PowerConnect, Force 10 Switches
- Reference Architecture
- Deployment Guide
- Dell Service and Support

cloudera

Revolutionary Cloud and Big Data Solutions

DELL

но, есть поставщики, предлагающие комплексные пакеты, что включают в себя и сам Hadoop, и утилиты по управлению им и автоматизации его развертывания на большое количество серверов. В частности, есть платформа Cloudera, есть дистрибутив Hadoop от Intel. Такие пакеты могут поставляться как бесплатно, так и на коммерческой основе с предоставлением технической поддержки. На этом уровне стэка “Big Data” компания Dell предлагает решения в сотрудничестве с компанией Cloudera, что позволяет заказчикам быстрее и проще реализовать аналитическую платформу, получить поддержку производителя и снизить риски при внедрении.

Четвертый уровень предложений – более сложный, это уже **аналитические базы данных**, создающиеся поверх хранилища данных. Вышеупомянутые пакеты программных разработок от сторонних поставщиков зачастую включают в себя не только платформу, но и саму базу данных. К этому же классу можно отнести еще несколько так называемых нереляционных баз данных NoSQL, например разработку MongoDB. Она может использоваться как хранилище для большого количества неструктурируемого контента.

В принципе, все, о чем говорилось выше, – это в основном программные разработки, которые позволяют так или иначе хранить данные в большом объеме. Однако производительность определяется не только кодом платформы, но и наличием аппаратных средств. Сколько необходимо для обработки данных серверов, сколько дисков? Как правильно подобрать соотношение емкости к процессорной мощности? Для упрощения решения этой задачи разработаны готовые **программно-аппаратные (Appliance) решения**. Например, на рынке хорошо известны такие решения, как Oracle Exadata, Greenplum Appliance, Teradata. Создавая решение Big Data на основе этих продуктов, на выходе заказчики получают готовую аналитическую машину, которая включает в себя как аппаратное, так и программное обеспечение. При этом аппаратная часть этих решений мало чем отличается от того, что было рассмотрено на первом

уровне предложений. Она может представлять собой обычные диски, находящиеся либо в составе вычислительных узлов и подключенные по внутренней шине, либо же это внешняя полка с дисками, подключаемая напрямую к вычислительному узлу. Несмотря на кажущееся отсутствие инновационности ценность Appliance-решений заключается в том, что количество оборудования, типы дисков, мощность процессоров, соотношение процессоров и емкости заранее рассчитано производителем и оптимизировано под оптимальную работу ПО. Соответственно, пределы, шаг и способ масштабирования в этом случае определены и контролируются производителем, что позволяет заказчику уже не думать на эту тему. Такой подход экономит время и снижает риски при внедрении.

При этом нельзя забывать, что на выходе, вне зависимости от того, какой подход был выбран, заказчик получит всего лишь машину, которая позволит работать с большим объемом данных. Сама по себе эта машина не будет выполнять те задачи, которые заказчик намеревается выполнять с ее помощью. Именно поэтому следует выделять еще один уровень предложений. Этим уровнем являются люди – **квалифицированные специалисты**, которые умеют анализировать, извлекать и придумывать алгоритмы для работы с Большими Данными. Собрав или купив платформу, заказчик получает лишь основу, на базе которой специалисты будут разрабатывать различные специфичные алгоритмы аналитики, релевантные каждой конкретной задаче, будь то маркетинг, геологоразведка, онлайн-трейдинг, финансовые услуги и т.д. В идеале под каждую компанию, под каждую модель бизнеса создается свой собственный алгоритм.

На данный момент в отрасли катастрофически не хватает таких специалистов по аналитике Больших Данных. Тем не менее, спрос порождает предложение, и ИТ-индустрия начинает довольно быстро растить соответствующие кадры. В России уже появляются соответствующие образовательные программы, связанные с подготовкой специалистов данного класса. Многие российские компании прилагают целенаправленные усилия в этой области и усиленными темпами готовят “шестой уровень” во всей этой структуре, без которой вся отрасль больших вычислений не имеет смысла.

Так или иначе, подводя итог, можно еще раз подчеркнуть, что в настоящее время на рынке имеется как минимум пять классов предложений в области технологий для хранения и аналитики Больших Данных, и в каждом классе предлагается достаточное количество решений от разных производителей.

**Михаил Орленко, руководитель  
департамента корпоративных решений,  
компания Dell в России,  
Казахстане и Центральной Азии**