大数据分析系统的建设

ruyue.ma@gmail.com

推销两个观点

- 数据系统
- 分层设计
- 为什么推销这两个观点?

数据系统

What is a data system? A system that manages the storage and querying of data.

Query = Function(All data).

Sometimes you retrieve what you stored Oftentimes you do transformations, aggregations, etc.

MapReduce is a framework for computing arbitrary functions on arbitrary data.

@nathanmarz, the author of Storm

Mysql是数据系统,包括查询层(SQL)和存储层。

HBase、HDFS、 NoSQL都至多算是存储层。

cloudera

Cloudera Impala: HyperTable

Cloudera Impala: How is Hypertable compared to Impala?

Hypertable: http://www.hypertable.com 4

2 Answers



Doug Judd



17 votes by Mark Johnson, Rebecca Ritter, Christoph Rupp, (more)



Hypertable is a direct technology competitor to Impala. The goal of the Hypertable project, from its inception, has been to build a scalable database to store, analyize, and serve massive datasets. As part of this effort, we are building all of the powerful query capabilities that you find in a traditional relational database, which means full SQL support. Todd was misinformed when he made his response. The job of Hypertable is **not** to act as a simple storage component, but to be a full-featured, high performance, scalable database.

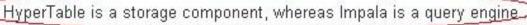
2+ Comments - 1:54 on Tue Nov 06 2012



Todd Lipcon, Employee since Feb '09



11 votes by Denis Arnaud, Patrick Angeles, Alon Amit, (more)





分层设计

- 垂直分层
 - 。 时效性库, 小时库, 天级库, 周库
- 水平分层
 - OLTP, OLAP
 - SQL, NoSQL
- 为什么分层?
 - 。减少设计复杂性
 - 。减少使用运维复杂性
 - 。资源效率使用最高

为什么推销这两个观点

- 小数据->大数据
- 大数据系统也是storage+query
- 大数据需要分层考虑

输出数据

实时计算层(Storm)

实时存储层(HBase)

批量计算层(MapReduce)

批量存储层(HDFS)

输入数据

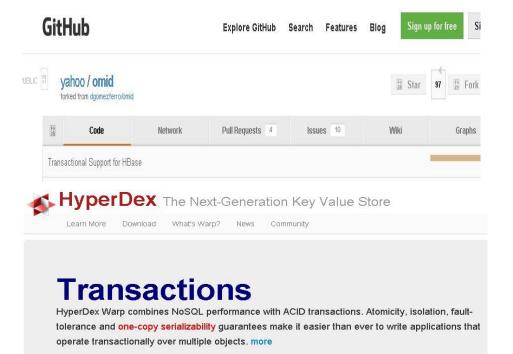
最近两个趋势

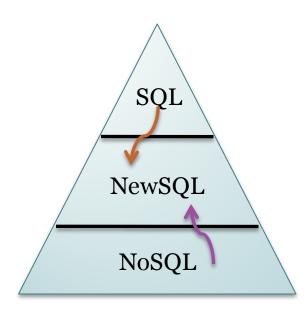
- NewSQL
- Interactive Analysis
- 说明了什么

NewSQL

• NoSQL太过原始, SQL容量性能有限

Megastore: Providing Scalable, Highly Available Storage for Interactive Services (Google 2011) F1 - The Fault-Tolerant Distributed RDBMS Supporting Google's Ad Business (Google 2012)





Interactive Analysis

• Hive响应太慢,数据库容量性能有限

Dremel: Interactive Analysis of Web-Scale Datasets (Google 2010)

Tenzing: A SQL Implementation On The MapReduce Framework (Google 2011)

PowerDrill: Processing a Trillion Cells per Mouse Click (Google 2012)





Enterprise Hadoop

Products

Hadoop Training

ommi

The Stinger Initiative: Making Apache Hive 100 Times Faster



February 20th, 2013

Alan Gates





WHY CLOUDERA

PRODUCTS

SOLUTIONS

PARTNERS

RESOURCES

SUPPORT

ABOUT

Hadoop & Big Data

Our Cuetomore

Cloudera Impala: Real-Time Queries in Apache Hadoop, For Real

by Marcel Kornacker & Justin Erickson | October 24, 2012 | # 14 comments | Tweet

Apache Drill Distributed system for interactive analysis.

Apache Drill (incubating) is a distributed system for interactive analysis of large-scale datasets, based on Google's Dremel. Its goal is to efficiently process nested data. It is a design goal to scale to 10,000 servers or more and to be able to process petabyes of data and trillions of records in seconds.

MemSQL, The Real-Time Analytics Platform.

MemSQL's real-time analytics platform is built on the world's fastest, most scalable in-memory database, capable of simultaneously handling real-time transactions and analytic workloads. MemSQL unleashes the full potential of Big Data by consuming and returning data instantly.

Making Data Work



Shark: Real-time queries and analytics for big data

Listen

Print

Shark is 100X faster than Hive for SQL, and 100X faster than Hadoop for machine-learning

by Ben Lorica | @bigdata | Comment | November 27, 2012

说明了什么

- 大数据的开源社区正在向数据库厂商发起挑战
- OLTP: 难度会稍大,撬动更多的是mysql、postgresql的领地
- OLAP: 很有希望
 - 。成本昂贵
 - 。稳定性要求低
 - 。数据量大
 - 。时效性低
 - 。不是不可缺少的组件

今天的重点:

大数据分析

热词榜

Amazon RedShift Teradata

EMC Greenplum
IBM Netezza HP Vertica
SAP Hana Oracle Exadata

Tajo EMC Hawk

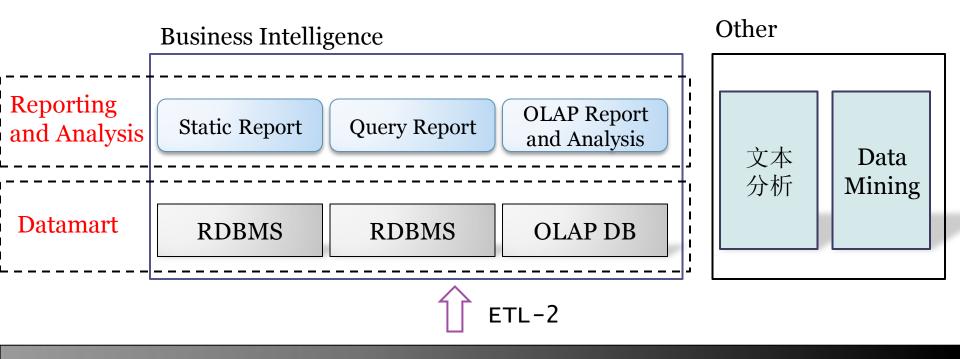
Stado citusdata Impala

Stinger/Tenz Pig/Hive

HPCC System Salesforce Phoenix

大数据分析架构

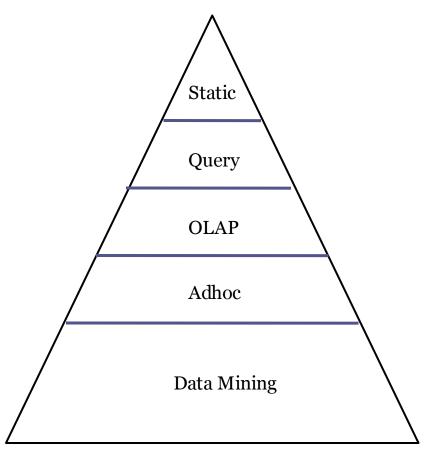




Dataware House



大数据分析发展趋势



从上往下:

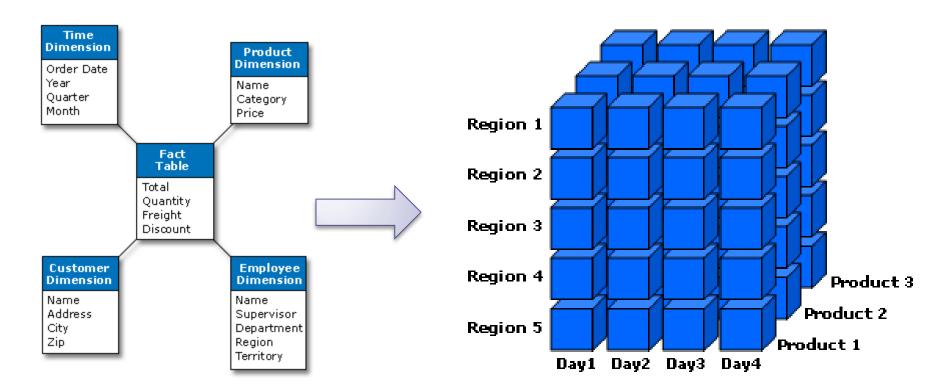
- 1. 数据量越来越大,维度越来越多
- 2. 交互性越来越难做
- 3. 技术难度越来越大
- 4. 以人为主->以机器为主
- 5. 用户专业程度越来越高,越来越少

非结构化数据

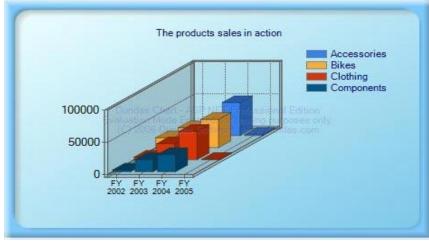
- 非结构化数据进行结构化后,利用原有技术分析
- 直接文本分析
 - □ 百度热搜词 static report
 - □ 用户query分析 query report
 - 。搜索引擎 OLAP多维分析
 - □ MapReduce上的调研作业 adhoc
 - □ 新闻聚类 data mining

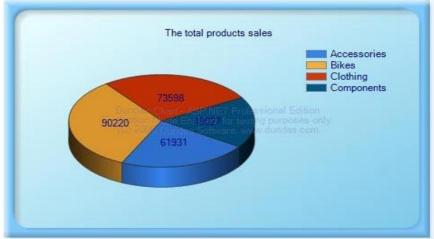
OLAP技术难点

- 多维分析: rollup, drill-down, slicing和dicing
- 各类维度组合,并提供交互式响应





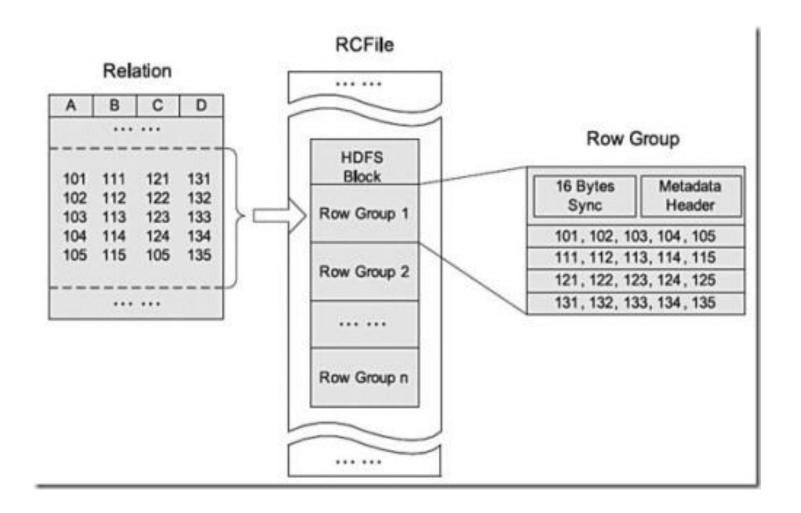




OLAP技术难点 - 解决手段

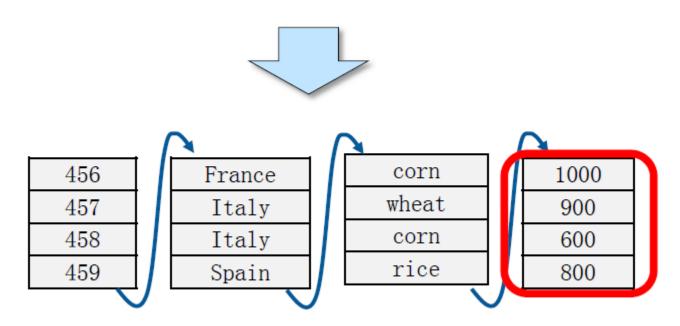
- 减少不必要的列读写
 - 。行列混合
 - 。列式存储
- 减少不必要的行读写
 - □ hyperdex 多维hash
 - infobright knowledge grid
- 压缩
- 预先计算
 - 。块级别的
 - 。物化视图

减少不必要的列读写 - 行列混合



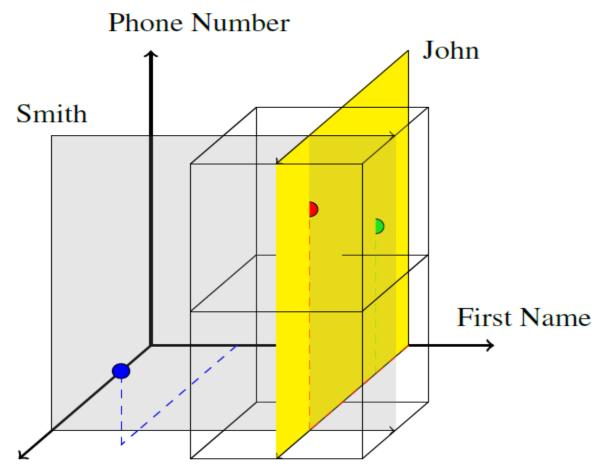
减少不必要的列读写 - 列式存储

Order	Country	Product	Sales
456	France	corn	1000
457	Italy	wheat	900
458	Italy	corn	600
459	Spain	rice	800



Column order organization

减少不必要的行读写 – 多维hash



Last Name

减少不必要的行读写 – infobright knowledge grid

Knowledge Grid

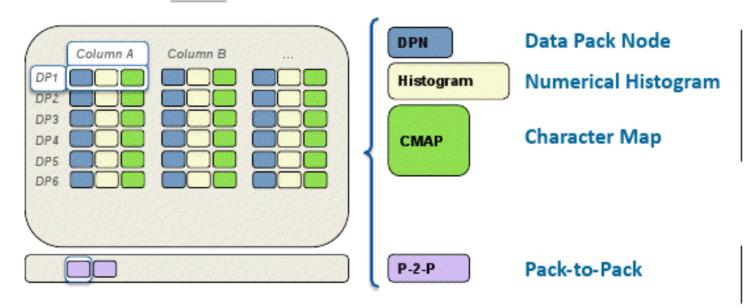
ŝ

applies to the whole table

Knowledge Nodes

built for each Data Pack

Information about the data



Built during LOAD

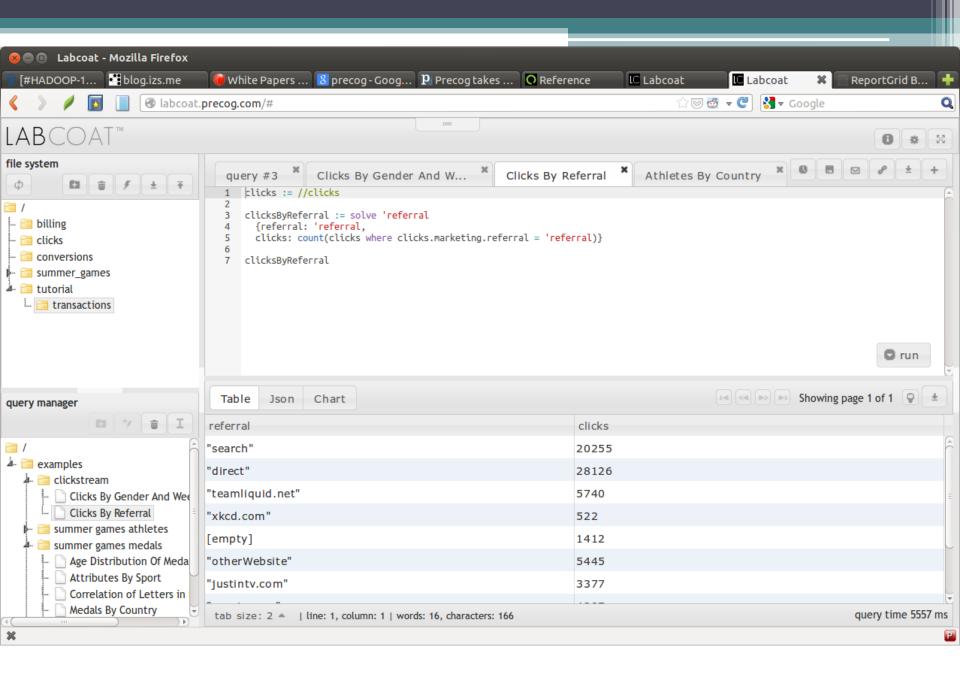
Built using JOIN

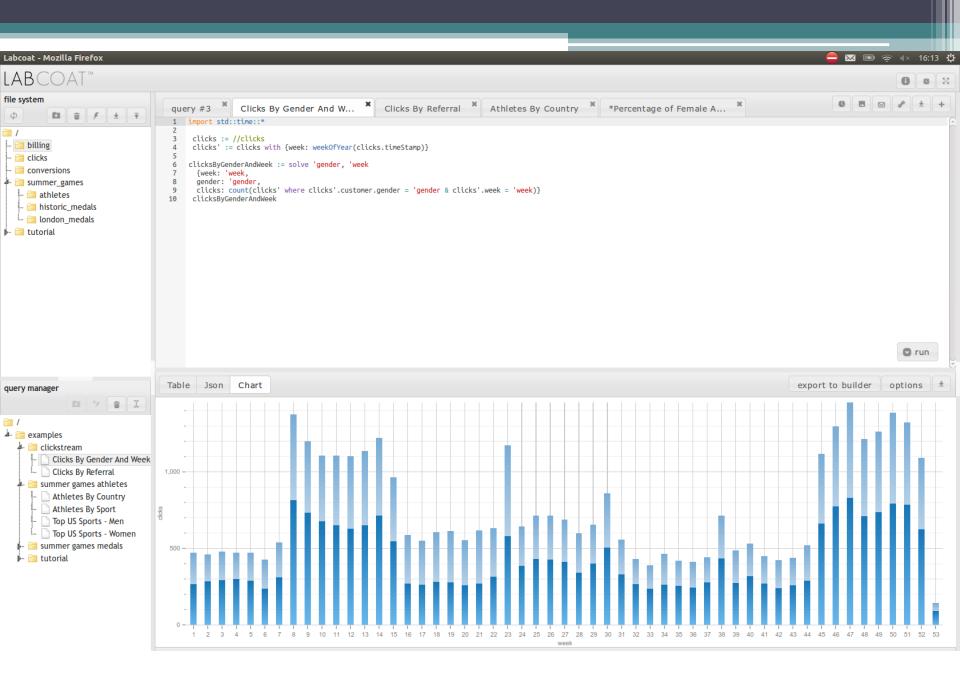
预先计算

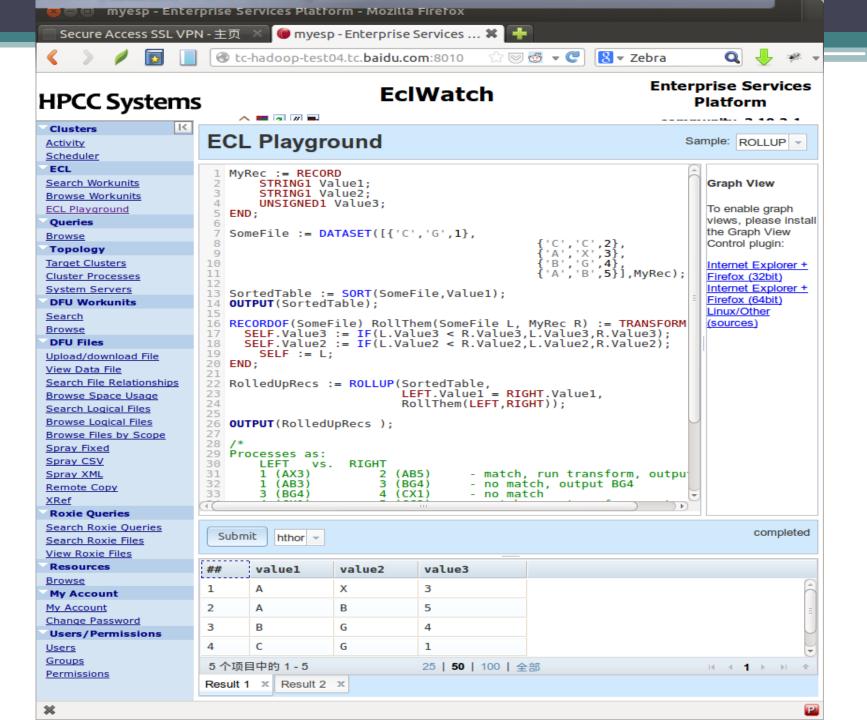
- 块级别的
 - 。对每一个数据块,提前计算好其max,min,sum,count等。
- 物化视图
 - 。提前计算好需要的几个维度的rollup表

Adhoc技术难点

- 任意维度分析: 存储优化, 等同OLAP
- 交互式响应
 - MapReduce太慢
 - Impala
- · 任意分析: 简单的SQL可能并不好用
- 方便的查询分析编写环境和展现工具
 - 。数据可能导入OLAP做进一步分析





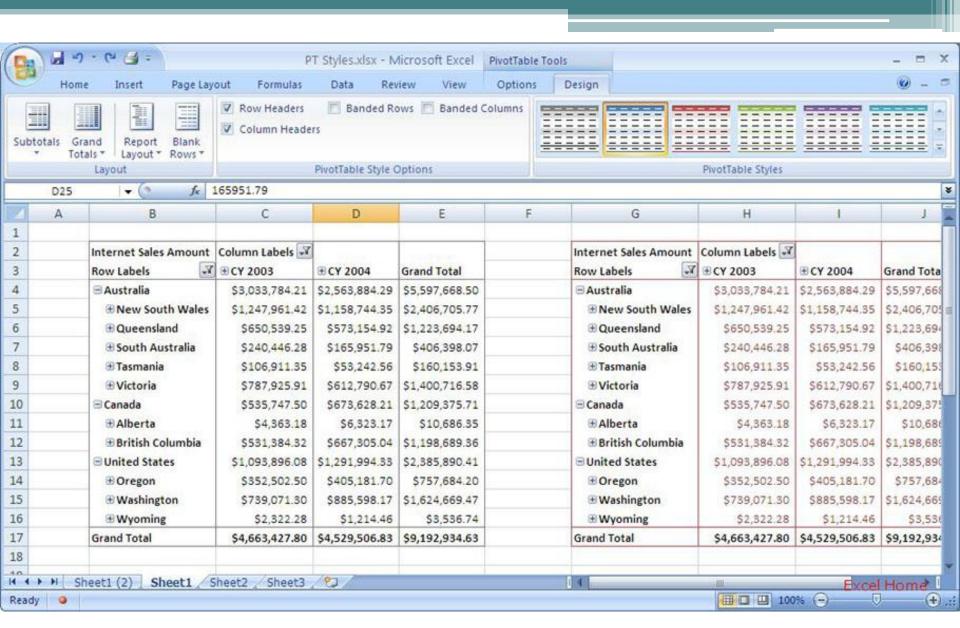


数据分析系统搭建 - 小系统

MS Excel (BI)

PentahoBI/SpagoBI

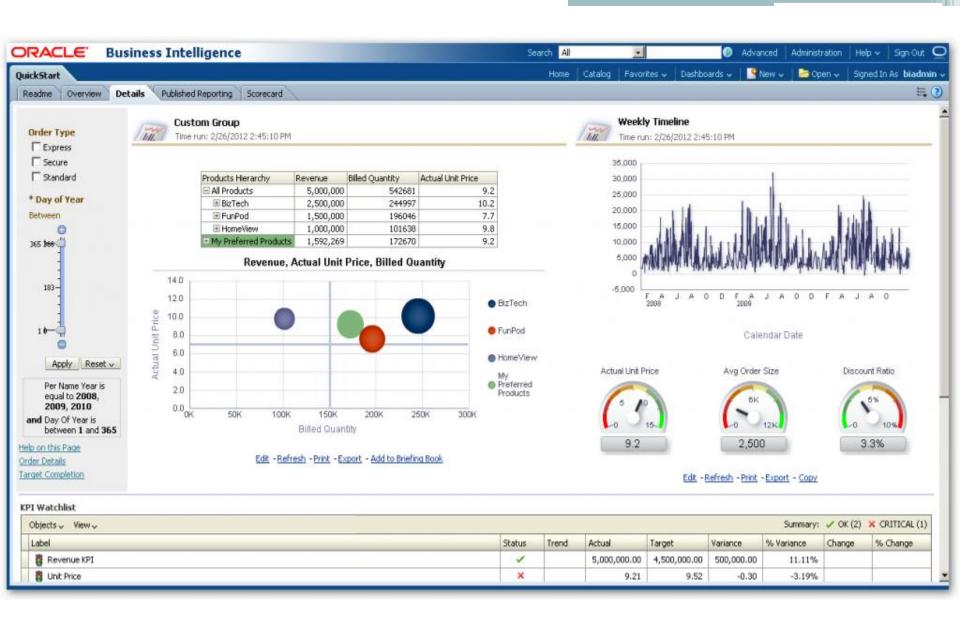
MySQL/Postgres/Infobright





大数据分析系统搭建 - 商业版

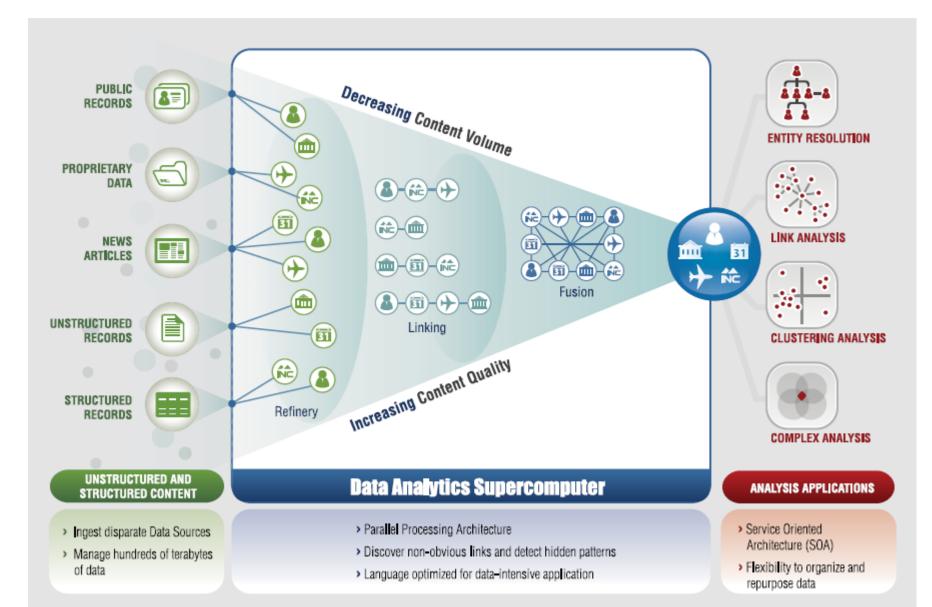
- Oracle BIEE + Oracle Exadata
- 其它产品
 - Greenplum
 - SAP HANA
 - Netteza



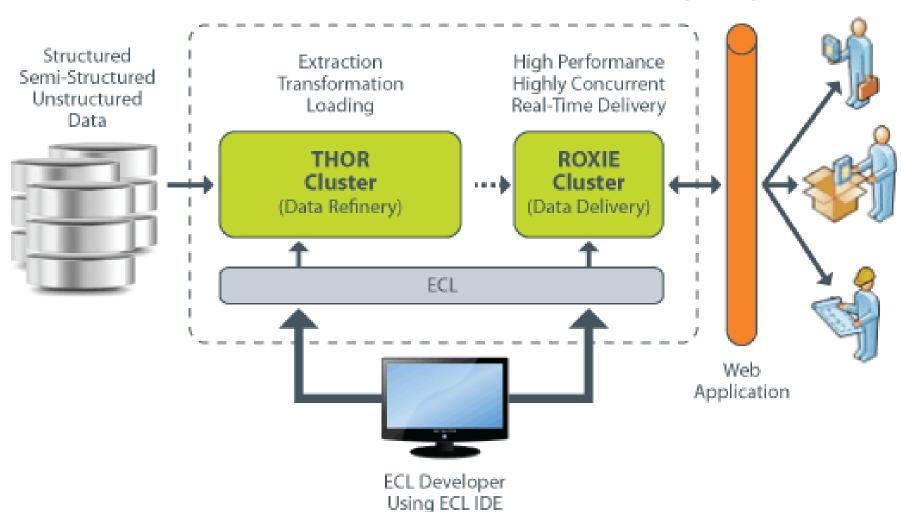
大数据分析系统搭建 - 开源版

- 开源还没有很成熟的产品来构建大数据下的 OLAP
- 短期解决
 - □商业产品
 - 。交互性强,访问量大:转为查询请求放入SQL或 NoSQL中查询
 - 。交互性要求不高,访问量少的: 转为利用 Hive/Impala来做

中等规模分析方案-HPCC Systems



HIGH PERFORMANCE COMPUTING CLUSTER (HPCC)

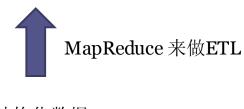


大规模分析方案 - Hadoop

Static Report/Query Report/小OLAP分析
SQL DB/NoSQL DB



HDFS





非结构化数据

结构化数据

谢

Q&A