

Big Data Infrastruktur – The Oracle Way

Daniel Steiger
Trivadis AG
Zürich-Glattbrugg

Schlüsselworte

Engineered Systems, Appliance, Big Data, Hadoop, Cloudera

Einleitung

Im Jahre 2006, also vor zehn Jahren, wurde der erste Release von Hadoop veröffentlicht. Zwei Jahre später wurde die Firma Cloudera gegründet – eine Unternehmung welche sehr früh auf die damals noch recht neue Big Data Technologie gesetzt hat. Weitere drei Jahre später hat Oracle sein Engineered Systems Portfolio um die Big Data Appliance erweitert und im darauffolgenden Jahr die Cloudera Hadoop-Distribution darauf integriert.

Damit war die Basis für eine Architektur gelegt, die heute unter dem Begriff „Oracle Big Data Management System Architektur“ zusammengefasst ist. In diesem Vortrag wird die Oracle Big Data Management System Architektur im Überblick vorgestellt, der Hardware- und Software-Stack der Big Data Appliance (BDA) erläutert und Erfahrungen zum Setup einer BDA weitergegeben. Abgerundet wird die Präsentation mit der Vorstellung eines konkreten Anwendungsfalles.

Oracle Big Data Management System Architektur

Die Oracle Big Data Management System Architektur (Abbildung 1) kombiniert die Performanz von Oracle's relationaler Datenbank und SQL-Engine mit dem kosteneffizienten und flexiblen Datenspeicher in Hadoop und NoSQL. Mit Hilfe der Softwareprodukte Oracle Big Data SQL, Oracle Golden Gate, u.a.m. steht eine integrierte und auf die Anforderungen für die Verwaltung von grossen und heterogenen Daten ausgerichtete Architektur zur Verfügung.

Die Oracle Big Data Appliance übernimmt dabei die Rolle des intelligenten Datenreservoirs. Damit können Daten aus unterschiedlichsten externen Quellen gesammelt, organisiert und ausgewertet werden. Mit Oracle Big Data Discovery werden Daten beispielsweise sehr effizient „an Ort“ ausgewertet und dargestellt, ohne dass Netzwerk-Latenzzeiten zu Verzögerungen führen. Für die performante Datenanalyse und –verarbeitung auf einem Oracle Exadata-System oder einer Oracle Exalytics In-Memory Machine wird die BDA über eine schnelle Infiniband-Verbindung angebunden.

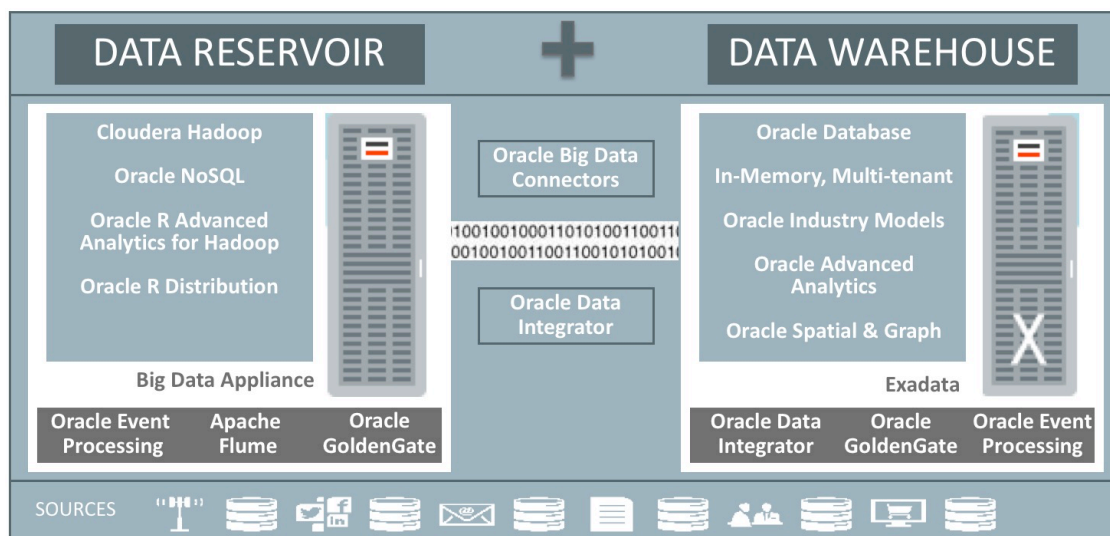


Abbildung 1: Oracle Big Data Management System Architecture (Quelle: Oracle)

Oracle Big Data Software

Der Oracle Big Data Software-Stack umfasst die für die Integration von traditionellen RDBMS-basierten Systemen und Big Data Komponenten wie Hadoop und NoSQL-Datenbanken notwendigen Softwareprodukte und Konnektoren. Darunter fallen folgenden Produkte:

- **Oracle Big Data SQL**
Oracle Big Data SQL erlaubt Abfrage mittels SQL auf Hadoop HDFS, NoSQL-Datenbanken und Oracle Datenbanken. Die SQL-Statements werden dabei so ausgeführt, dass die Optimierung auf dem Quellsystem – also beispielsweise auf einem Hadoop-Node - stattfindet. Innerhalb der (relationalen) Oracle Datenbank werden die Hadoop oder NoSQL Daten als externe Tabellen definiert. Mittels Nutzung eines Hive Metastores kann Oracle dynamisch über eine mögliche Parallelisierung und über die Lese-Semantik entscheiden. Mit dem Oracle Big Data SQL Agent sind zudem Features wie Smart Scan for Hadoop, Storage Indexes und Local filtering and Caching auch auf einem Non-Oracle Engineered System verfügbar:
- **Oracle Big Data Discovery**
Oracle Big Data Discovery wird auch als "Visual Face of Big Data" bezeichnet. Big Data Discovery nutzt Apache Spark – ein hochperformantes Cluster Computing Framework – um sehr grosse Datenmengen zu verarbeiten.
- **Oracle Data Integrator for Big Data (ODI)**
ODI for Big Data dient der Datentransformation und der Informationsanreicherung innerhalb des Datareservoirs. ODI generiert dabei native Code der auf der Hadoop-Plattform ausgeführt wird. Dies erlaubt Business -und Daten-Verknüpfungen herzustellen ohne dass eine spezifische Programmiersprache wie HiveSQL, Pig Latin oder Map Reduce gelernt werden muss.
- **Oracle GoldenGate for Big Data**
GoldenGate for Big Data ist ein bewährtes und flexibles Werkzeug um Daten von einem traditionellen Quellsystem „minimal-invasiv“ zu einem Big Data Zielsystem zu liefern. Mit dem Release 12.2 von GoldenGate for Big Data werden diverse neue Protokolle wie JSON, AVRO und XML, sowie Native Java Replication und Kafka unterstützt. Es eignet sich auch für den Einsatz für Real-Time Streaming Analytics

Die erwähnten Softwarekomponenten, aber auch die Komponenten welche auf der Big Data Appliance mit dem Cloudera Distribution ausgeliefert werden, arbeiten nach dem Prinzip der „Data Locality“. Dies bedeutet, dass die Verarbeitung der Daten lokal, d.h. am Ort der Datenhaltung erfolgt und dementsprechend effizient ist.

Oracle Big Data Appliance

Die Big Data Appliance (BDA) ist die Schlüsselkomponente in der von Oracle propagierten "Oracle Big Data Management System"-Architektur. Die Big Data Appliance ist eine Hardware- und Software-Komplettlösung für Big Data Anforderungen. Sie steht in drei Konfigurationen zur Verfügung:

- Das Starter Rack besteht aus 6 Compute/Storage-Knoten
- Ein Full Rack ist mit 18 Compute/Storage-Knoten bestückt
- In einer Multi-Rack-Konfiguration können bis maximal 18 Racks zu einem Gesamtsystem miteinander verbunden werden

Der X6-2 Compute/Storage-Knoten ist der eigentliche Building-Block der Big Data Appliance. Er verfügt über folgende Spezifikation:

- 2 x 22-Core (2.2GHz) Intel ® Xeon ® E5-2699 v4
- 8 x 32GB DDR4-2400 Memory (expandable to maximum of 768GB per node)
- 12 x 8TB 7,200 RPM High Capacity SAS Drives
- 2 x QDR 40Gb/sec InfiniBand Ports
- 4 x 10 Gb Ethernet Ports
- 1 x ILOM Ethernet Port

Die Compute/Storage-Knoten sind über den InfiniBand Network Backbone miteinander verbunden. InfiniBand garantiert sowohl eine sehr schnelle Verbindung zwischen den einzelnen Cluster-Knoten, wie auch zwischen mehreren BDAs und weiteren Oracle Engineered Systems wie der Exadata, Exalogic, SPARC SuperCluster und Exalytics.

Oracle Big Data Appliance Softwarestack

Mit der BDA wird die Cloudera Enterprise Data Hub Edition ausgeliefert. Darin enthalten ist folgende Software:

- Apache Hadoop (CDH)
- Cloudera Impala
- Cloudera Search (Apache Solr)
- Apache HBase and Apache Accumulo
- Apache Spark
- Apache Kafka
- Cloudera Manager

Setup

Das Setup einer Big Data Appliance beginnt mit der physischen Installation des Racks in einem geeigneten Serverraum. Der Safety and Compliance Guide, sowie die Site Checklist leisten dafür gute Dienste.

Sobald die Appliance am Strom und am Netzwerk angeschlossen ist, kann mit dem Mammoth-Utility das eigentliche Software-Deployment durchgeführt werden:

```
cd /opt/oracle/BDAMammoth
mammoth -s 1 cdh
```

Tipps

Für die Installation, Patching und den Unterhalt einer Big Data Appliance sind sehr gute technische Kenntnisse wie Linux-Administration, ssh, X-Server, scp, Netzwerk, Storage, etc. eine wichtige Voraussetzung. Erfahrung mit anderen Oracle Engineered Systems, wie Exadata oder der Database Appliance, sind ebenfalls sehr hilfreich, weil viele Aufgaben, wie beispielsweise Patching und SW-Deployment, sehr ähnlich sind.

Im Weiteren haben wir die Erfahrung gemacht, dass die Integration der BDA mit Exadata über Infiniband ihre Tücken haben kann. Wir empfehlen auf jeden Fall die entsprechenden My Oracle Support Seiten (siehe Referenzen am Schluss des Dokuments) zu konsultieren.

Use Case

Ein Kunde aus dem Sportbereich stand im Rahmen eines Projektes zur Entdeckung von Betrugserkennung im Spielwettbewerb vor der Herausforderung, Daten aus unterschiedlichsten Quellen zu analysieren und Auffälligkeiten bezüglich der Leistung einzelner Spieler zeitnah zu erkennen und darauf reagieren zu können.

Da beim Kunden bereits Exadata-Systeme im Einsatz sind, hat sich er sich entschieden auch für das Big Data System eine Oracle-Lösung einzusetzen. Massgebliche Argumente für die Wahl der Oracle Big Data Appliance waren folgende:

- Integrationsvorteile, weil alle System vom gleichen Hersteller sind
- Engineeringaufwand praktisch null
- „Fast start to Big Data" dank umfassenden Software-Paket (Cloudera Distribution und Oracle-Software) für die Datenanalyse und die Integration in bestehende Systeme
- Kostengünstige Minimalkonfiguration, welche bei Bedarf skaliert werden kann

Zusammenfassung

Mit der Oracle Big Data Management System Architektur werden alle Phasen der (Big Data) Datenverarbeitung – von der Erfassung, Bereinigung, Ablage, über Real-Time-Analyse, Event-Processing, Visualisierung, bis zur Langzeit-Speicherung abgedeckt.

Als Data Reservoir erfüllt die Big Data Appliance die Rolle als multifunktionales Buildingblock in der Aufnahme und Verarbeitung von grossen, heterogenen Datenmengen. Die Hardware der BDA erfüllt dabei die technischen und funktionalen Anforderungen einer Enterprise-Infrastruktur und ist mit einem sehr umfangreichen Softwarepaket (Cloudera Enterprise Data Hub Edition) ausgestattet.

Die Big Data Appliance wird sich mit zukünftigen Releases – vor allem im Softwarebereich - als sog. Converged System weiterentwickeln und verbessern. Man darf dabei nicht vergessen, dass die Basistechnology Hadoop im Vergleich zu etablierten Technologien wie ein RDBMS noch recht jung ist. Last but not least: es ist sicher keine falsche Strategie bei einem Big Data Projekt auf das Zusammenspiel von Engineered System Expertise und Data Analytics Knowhow zu setzen.

Referenzen und weiterführende Informationen

- Information Center: Oracle Big Data Appliance, My Oracle Support Doc ID 1445762.2
- Oracle Big Data Appliance Installation Frequently Asked Questions, My Oracle Support Doc ID 1518939.1
- Oracle Big Data Appliance Patch Set Master Note, My Oracle Support Doc ID 1485745.1
- Oracle Big Data Appliance Documentation,
<http://www.oracle.com/technetwork/database/bigdata-appliance/documentation/index.html>
- An Enterprise Architect's Guide to Big Data – Reference Architecture Overview,
<http://www.oracle.com/technetwork/topics/entarch/oracle-wp-big-data-refarch-2019930.pdf>

Kontaktadresse:

Daniel Steiger
Trivadis AG
Sägereistrasse 29
CH-8152 Glattbrugg

Telefon: +41 58 459 50 88
E-Mail: daniel.steiger@trivadis.com
Internet: www.trivadis.com