

Data Warehouse

Datenstrukturen

- Stern/ Star Schema
 - Es steht nicht die Normalisierung sondern die Optimierung auf effiziente Leseoperationen
 - Im Zentrum die Faktentabelle und aussen rum sind die Dimension
 - Sternschem ist denormalisiert, Anomalien und erhöhter Speicherbedarf ist möglich
- Snowflake
 - Hier müssen allerdings mehrstufige Dimensionstabellen über Join-Abfragen verknüpft werden.
 - Hat die Dimensionen in der 3.NF
 - Erweiterte Unterteilungen etc.
- Galaxy
 - Teilt die gleichen Dimensionen mit mehreren Faktentabellen

Referenzmodell in der Praxis umgesetzt

- Hub and Spoke Architektur
 - Nabe Speicher Architektur
 - Zentrales DWH (HUB) mit angeschlossenen DM (Spokes)
- Legacy / Middleware
 - aus verschiedenen Legacy Systemen direkt ins Data Warehouse, daraus in die DataMarts via Middleware in Anwendungen (Analysen)
- High Level System Architecture
 - Alles auf ändlicher Ebene nahe am Präsentations Layer

Dies soll unzeigen wie in wir ein DWH in eine Unternehmen integrieren können.

Bestandteile des Referenzmodells Definitionen

- Reportingbereich SSRS
 - Im Referenzmodell nicht erwähnt
 - SQL Server Reporting Services SSRS
- Auswertungsbereich SSAS
 - SQL Server Analysis Services = SSAS
 - Ableitungsdatenbank
 - Zentrales Data DWH
 - möglichst homogener Datenbestand
 - Data Mart
 - Teilkopien von DW für u.a Abteilungen
 - Cube
 - für mehrdimensionale Abfragen
- Integrationsbereich SSIS
 - SQL Server Integration Service = SSIS
 - periodische Transformation od. Homogenisierung
 - Arbeitsbereich = Staging Area
- Verwaltungsbereich
 - DB's
 - Flatfiles
 - Webserver
 - CRM- Systeme
 - Repository
- Integrationsbereich SSIS
 - periodische Transformation od. Homogenisierung
 - Arbeitsbereich = Staging Area
- Verwaltungsbereich
 - DB's
 - Flatfiles
 - Webserver
 - CRM- Systeme
 - Repository

Bestandteile der Referenzarchitektur im Detail

- 1. Datenquellen - Externe Daten
 - Daten aus Legacy Systemen
 - Daten aus Anwendungsprogrammen SAP
 - Zeichencodierung ASCII, ANSI, UTF 8
 - Trennzeichen Abstand, Tabulator, Semikolon
 - Zeilenschaltung (Paragraph) CR/ LF Windows, LF Unix+Mac
 - Meist Falsche Tabellen oder Flache Files
 - XML Datenströme soll und muss heute Standard werden
- 2. Arbeitsbereich Staging Area
 - Integriert Daten, die Heterogenen
 - Das E von ETL
 - Date werden mit Hilfe der Metadaten zusammengeführt und abgelegt
 - Metadaten Respository
 - Datenbank Tabellen, Daten Stammen aus sehr unterschiedlichen Systemen. Alle Notwendigen Beschreibungen zum System können zur Umwelt darin enthalten sein.
 - Arbeitsdepot, das flüchtige Zwischendepot
 - Auch für Simulationen, Berechnungen usw
 - Zwischenstufe von Quellsystem zum DWH
 - Auch als Last und Leistungsausgleich
 - z.b für voranalysen
 - Data Cleaning & Data Wrangling
 - Daten in der BasisDatenbank sind Bereinigt
 - ODS und Stating Area kann und muss nicht immer Zusammen sein
 - Integrationsbereich
- 3. Basis DB (Operational Data Sorage ODS)
 - Integrationsbereich
- 4. Ableitungsdatenbank Data Warehouse
 - Data Vault
- 5. Auswertungsbereich (Data Mart)
 - Teilkopie des DWH
 - bessere Leistung
 - Verlagerung der Rechenleistung
 - Weniger Network Traffic
 - Anwendungsorientiert nicht wie dwh das anwendungsneutral ist
 - Verwendung
 - Abteilung vs. Unternehmen
 - Anzahl
 - Ein Data Mart im Unterscheid zum Data Lake hat viel Potenzial. Man Stelle sich vor das man ein Wein auch in einem Delikatessenladen kauft für einen Speziellen Anlass. Der Spezialist im Laden kennt sich viel besser mit den Produkten aus als in einer Warenkette. Es kann sicher möglich sein beide Produkte an beiden Orten zu bekommen, aber im Delikatessenladen wird man höchstwahrscheinlich nicht sie schalchten weinse bekommen die man höchstwahrscheinlich kaufen würde in einem Supermarkt.
 - Unterschied zu DWH

DWH soll: "Reich an Daten, arm an Informationen und Wissen" wiederlegen

Herausforderungen

- DWH führ strukturierte und unstrukturierte Daten eine einem Ort zusammen
- Kein vorgegebenes Schema
- Zugriff abhängig von Semantik (Interprationsspielraum)
- hohes Datenvolumen, kein manuelles Extrahieren möglich
- Prioritisierung je nach Semantik und Kontext
- Datenqualität kann nicht Autonom gemessen werden

Unstrukturierte Daten

Schritte für logische und Physische Integration

- Extraktion
 - komplex
 - Information Retrieval
- Analyse für Integration (NLP - Natural Language Processing)
 - Kategorisieren, Auswertung von Dok./Flatfile
- Datenbereinigung (Data Cleansing)
 - Standard Cleaning
 - NLP
- Laden in DBMS, DWH
- ETL Extract Transform Load
 - Connect & Collect, Tmasform & Enrich, Publish, Monitor à la definition Microsoft

Definiton

Ein Datawarehouse ist eine konsistente, speziell analytische entwickelte Datenbank, welche ein Unternehmen massgeblich bei der Entscheidung hinsichtlich auf das Decision Mangement unterstützen kann. Die Hauptfunktion ist eine einheitliche Umgebung zu schaffen aus welcher BI Tools, sprich OLAP do. Multidimensionale Abfragen generiert werden können

Ein Data - Warehouse integriert aus vielen unterschiedlichen Quellen in einer für die Entscheidungsfindung optimierten Datenbank

