

DATA WAREHOUSE AUTOMATION MIT DATA VAULT

...bezeichnet eine Modellierungstechnik für das Data Warehouse

...folgt klaren Regeln

...hohe Flexibilität im Hinblick die Erweiterbarkeit

...legt den Fokus darauf, den stetig neuen Geschäftsanforderungen gerecht zu werden

...vereint Aspekte der rationalen Datenbankmodellierung mit der dritten Normalform sowie dem Sternschema

...vollständige Historisierung der Daten

Die Trennung der Konsolidierung und ihrer Business-Logik von der Datenintegration und Historisierung kann ein Teil der Verarbeitung standardisiert und damit automatisiert werden

Quellen: [1, 2]

WEITERENTWICKLUNG ZUM DATA VAULT 2.0

- Die ersten Veröffentlichungen zum Thema Data Vault gehen auf Linstedt um das Jahr 2000 zurück.
- Es folgte eine stetige Weiterentwicklung des Datenmodells, in dessen Zusammenhang auch der Begriff **Data Vault 2.0** entstand.
- Diese Erweiterungen wurden von Linstedt in einem Data Vault 2.0 Framework zusammengefasst.
- Es umfasst neben der Modellierung auch die agile Vorgehensweise bei der Implementierung und betrachtet die Architektur.
- Dabei werden zudem Performanceaspekte und unter anderem NoSQL Techniken bis hin zur Automatisierung vertieft.

Quellen: [2, 3, 8]

MODELLIERUNG

Die Grundidee des Data Vault folgt der Aufteilung der Daten in die Strukturelemente HUB, LINK und SATellit.

**HUB (Geschäftsobjekt)**

- Kernobjekt der Geschäftslogik
- Identifiziert die fachliche Identität
- Speichert:
  - Fachlichen Schlüssel (Business Keys)
  - künstlich generierten Primärschlüssel
  - Quellsystem und Ladezeitpunkt
- Speichert keine beschreibenden Attribute
- Zu beachten: direkte Verbindung zwischen Hubs sind nicht zulässig

**LINK (Geschäftsbeziehung)**

- Stellt die Beziehungen zwischen den HUBs dar
- Kann beliebig viele HUBs miteinander verknüpfen
- Kann keinen, einen oder mehrere SATelliten verknüpfen

**SATelliten**

- Informationen, die HUB und LINK beschreiben
- Speichert:
  - Beschreibende Attribute
  - Quellsystem und Ladezeitpunkt
  - Fremdschlüssel zu genau einem HUB oder LINK

Quellen: [2, 4]

MODELLIERUNGSPROZESS

Legende

- HUB
- LINK
- SAT

Fachliche Entitäten definieren:  
**HUBs**

Beziehungen modellieren:  
**LINKs**

Beschreibende Attribute festlegen:  
**SATelliten**

Quellen: [3, 4, 6, 7]

ARCHITEKTUR

Innerhalb des **Staging Layers** werden die Rohdaten aus den verschiedenen Quellsystemen eingesammelt

Hier wird das Data Vault modelliert. Es setzt sich aus den Rohdaten (Raw Data Vault) und harmonisierten und transformierte Daten (Business Data Vault), welche auf Geschäftsregeln und Laufzeitinformationen basieren zusammen

Im Mart Layer werden die Daten in die Form eines Star-Schema und/oder eines anderen Modellierungsverfahren transformiert und gespeichert. Diese Form der Daten stellen die Informationsquelle für die Visualisierung mittels der BI Plattformen dar.

QUELLSYSTEME

ERP

CRM

...

ETL

STAGING LAYER

Raw Data Vault

Business Data Vault

MARTLAYER

PRÄSENTATION

META DATEN

Quellen: angelehnt an [2, 4, 5, 9]

VORTEILE VON DATA VAULT FÜR UNTERNEHMEN

**ENTWICKLUNGSZEIT REDUZIEREN**  
Mithilfe von Data Vault lässt sich die Entwicklungszeit reduzieren. Dabei liegen die Vorteile vor allem bei der Implementierung der Business-Anforderungen.

**RETURN ON INVESTMENT**  
Unternehmen erzielen einen höheren Retrun on Investment (ROI).

**SKALIERBARKEIT**  
Durch den Einsatz von Data Vault steigt die Skalierbarkeit des DWH.

**NACHVERFOLGBARKEIT**  
Unternehmen können alle Daten bis zu ihrem Quellsystem nachverfolgen, wenn der Data Vault Ansatz genutzt wird.

**INKREMENTELL AUSBAUBAR**  
Dank des iterativen und agilen Entwicklungszyklus lässt sich das DWH flexibel ausbauen und erweitern.

**HISTORISIERUNG UND TIME TRAVELING**  
Durch die dauerhafte Speicherung der unveränderten Rohdaten, entsteht der Vorteil des Time Traveling, also stichtagsbezogene Auswertungen aus der Vergangenheit ermöglicht.

Quelle: [3]

GOLDENE REGELN

- Eine Hub-Tabelle gibt seinen Primary Key grundsätzlich nach außen
- Beziehungen zwischen Hubs gibt es nur über Links
- Rekursive Strukturen gibt es nur über Links
- Eine Link-Tabelle hat mindestens zwei Foreign Keys-Beziehungen
- Eine Link-Struktur kann eine Surrogatschlüssel-Beziehung aufweisen
- Eine Link-Tabelle kann eine unbegrenzte Anzahl an Hubs verbinden
- Eine Beziehung zwischen Links ist erlaubt
- Eine Link-Tabelle kann eine Satellit-Tabelle haben
- Ein Satellit kann nur eine Elterntabelle haben
- Ein Satellit kann nicht irgendeine Fremdschlüsselbeziehungen besitzen, mit Ausnahme des Primärschlüssels zu der Elterntabelle (Hub oder Link).

WANN IST DATA VAULT UNINTERESSANT?

Kleine Menge an Anforderungen  
Kurze Zeiträume, Berichte zu bauen und zu liefern  
Daten nur einmalig nutzen  
Nur ein Quellsystem (oder sehr wenige)  
Nur einen Analysten im Unternehmen  
Keine notwendigen Audit-Anforderungen zur Datenhistorisierung  
Keine notwendige Integration mehrere Data Centers  
Keine *near real time* Anforderungen  
Keine umfangreichen Batch Data  
Keine externe Daten außerhalb der eigenen Kontrolle  
Keine Trendanalysen  
Kein permanentes Re-Engineering

CHANCEN UND RISIKEN

CHANCEN	RISIKEN
<b>INTEGRATION VON DATEN AUS VERSCHIEDENEN QUELLSYSTEMEN</b> Die Quelldaten werden unter Verwendung gemeinsamer Geschäftsschlüssel, die in Hubs gespeichert sind, integriert. Die erforderlichen Geschäftsattribute werden in separaten Satelliten pro Quellsystem gespeichert. Dies erleichtert die Zusammenführung der Informationen für weitere Berichte.	<b>BUSINESS-KENNTNISSE ERFORDERLICH</b> Um erfolgreich Modelle mit Data Vault erstellen zu können, ist es wichtig, die Geschäftskontexte zu verstehen. Andernfalls ist das Risiko im Data Vault hoch, dass nur die Quelldaten kopiert und historisiert werden.
<b>PARALLELES LADEN VON DATEN AUS VERSCHIEDENEN QUELLSYSTEMEN</b> Es gibt keine vordefinierte Ladereihenfolge, die Daten können unabhängig voneinander in den Datenspeicher geladen werden.	<b>DATA VAULT KENNTNISSE ERFORDERLICH</b> Die Grundprinzipien von Data Vault müssen dem gesamten Projektteam bekannt sein. Die Entwickler müssen die ETL-Muster für die verschiedenen Objekte verstehen.
<b>VOLLSTÄNDIGE HISTORISIERUNG ALLER ATTRIBUTE</b> Die Versionierung aller Attribute in den Satelliten ermöglicht die Rückverfolgbarkeit aller Änderungen in der Vergangenheit und die Extraktion der Daten zu einem bestimmten Zeitpunkt.	<b>GROSSE TABELLENANZAHL</b> Viele Modellerweiterungen können Datenmodelle mit einer hohen Anzahl von Tabellen erstellen (Hubs, Links und Satelliten). Entsprechend steigt auch die Anzahl der ETL-Prozesse.
<b>LEICHTE ERWEITERBARKEIT DES DATENMODELLS</b> Zusätzliche Entitäten oder Attribute, die für neue Anforderungen verwendet werden, werden als zusätzliche Tabellen im Datendepot implementiert. Bestehende Tabellen werden normalerweise nicht geändert. Dies hilft, Datenmigration zu vermeiden.	<b>KOMPLEXE EXTRAKTION AUS DEM DATENDEPOT</b> Während das Laden von Data-Vault-Tabellen sehr einfach ist, kann das Extrahieren der Informationen zum Laden von Data-Marts umfangreicher sein. Für eine gute Leistung können Hilfstabellen erforderlich sein.
<b>EINFACHE UND EINHEITLICHE ETL-MUSTER</b> Das Laden von Hubs, Links und Satelliten erfolgt nach einheitlichen Regeln, die immer gleich aufgebaut sind.	<b>ENTSPRECHENDE GESCHÄFTSSCHLÜSSEL ERFORDERLICH</b> Die Bestimmung der geeigneten Geschäftsschlüssel ist eine der größten Herausforderungen in der Datenverwaltung. Ungeeignete Schlüssel erschweren die Integration verschiedener Quellen und erhöhen die Komplexität beim Laden von Data Marts.

Quelle [5]

Quellen:

[1] Michael Müller, Aktuelle Themen in die unternehmensweite BI integrieren: Architekturen in Data Vault 2.0, 2015, [https://www.tdwi.eu/fileadmin/tdwi/ext\\_wissen/whitepaper/MID\\_mueller\\_BIS\\_DataVault\\_2015\\_AktuelleThemen.pdf](https://www.tdwi.eu/fileadmin/tdwi/ext_wissen/whitepaper/MID_mueller_BIS_DataVault_2015_AktuelleThemen.pdf)

[2] Il-novum, Mit Data Vault zu mehr Agilität im Data Warehouse, [http://sigs.de/TDWI\\_Wissen/il-novum/whitepaper\\_data-vault.pdf](http://sigs.de/TDWI_Wissen/il-novum/whitepaper_data-vault.pdf)

[3] Dirk Lerner, Data Vault Modellierung, TDWI Wissen, [http://www.sigs.de/tdwi/Infografik/TDWI\\_Infografik\\_DataVault\\_Modelling.pdf](http://www.sigs.de/tdwi/Infografik/TDWI_Infografik_DataVault_Modelling.pdf)

[4] Linstedt, Daniel, and Michael Olschmke. Building a scalable data warehouse with data vault 2.0. Morgan Kaufmann, 2015.

[5] Trivadis AG, Data Vault – The ultimate problem solver?, [https://www.tdwi.eu/fileadmin/tdwi/1.0\\_Wissen/White\\_Paper/Trivadis\\_biGENIUS\\_Data\\_Vault-The\\_Ultimate\\_Problem\\_Solver.pdf](https://www.tdwi.eu/fileadmin/tdwi/1.0_Wissen/White_Paper/Trivadis_biGENIUS_Data_Vault-The_Ultimate_Problem_Solver.pdf)

[6] Schneider, Dani, Modellierung agiler Data Warehouses mit Data Vault, 2015, [http://dani.schneiderkennei.ch/pdf/DOAG\\_Data\\_Vault\\_Modellierung\\_Vortrag.pdf](http://dani.schneiderkennei.ch/pdf/DOAG_Data_Vault_Modellierung_Vortrag.pdf) <https://bigenius.info>

[7] Hultgren, Hans, Data Vault Modeling Guide – Introductory Guide to Data Vault Modeling, 2015, [https://www.tdwi.eu/fileadmin/tdwi/ext\\_wissen/whitepaper/centinnium-data-vault-modeling-guide.pdf](https://www.tdwi.eu/fileadmin/tdwi/ext_wissen/whitepaper/centinnium-data-vault-modeling-guide.pdf)

[8] Cramer, Oliver; Lerner, Dirk, Neue Wege in der Datenmodellierung - Data Vault heißt die moderne Antwort, 2015, [https://www.tdwi.eu/fileadmin/tdwi/ext\\_wissen/whitepaper/ITGAIN\\_lerner\\_BIS\\_DataVault\\_2015.pdf](https://www.tdwi.eu/fileadmin/tdwi/ext_wissen/whitepaper/ITGAIN_lerner_BIS_DataVault_2015.pdf)

[9] Bauer, Lutz, Mehr Umsetzungs-Geschwindigkeit und Flexibilität für Ihr Data Warehouse Industrialisierter Data Vault, 2015, [https://www.tdwi.eu/fileadmin/tdwi/ext\\_wissen/whitepaper/MT\\_bauer\\_BIS\\_DataVault\\_2015.pdf](https://www.tdwi.eu/fileadmin/tdwi/ext_wissen/whitepaper/MT_bauer_BIS_DataVault_2015.pdf).

tdwi.eu

Autor:

Prof. Carsten Feldten  
Technische Universität Bergakademie Freiberg (Sachsen)