

dani SCHNIDER  
claus JORDAN  
peter WELKER  
joachim WEHNER

# DATA WAREHOUSE BLUEPRINTS

Business Intelligence in der Praxis



HANSER

## Bleiben Sie auf dem Laufenden!



Unser **Computerbuch-Newsletter** informiert Sie monatlich über neue Bücher und Termine. Profitieren Sie auch von Gewinnspielen und exklusiven Leseproben. Gleich anmelden unter



[www.hanser-fachbuch.de/newsletter](http://www.hanser-fachbuch.de/newsletter)



**Hanser Update** ist der IT-Blog des Hanser Verlags mit Beiträgen und Praxistipps von unseren Autoren rund um die Themen Online Marketing, Webentwicklung, Programmierung, Softwareentwicklung sowie IT- und Projektmanagement. Lesen Sie mit und abonnieren Sie unsere News unter



[www.hanser-fachbuch.de/update](http://www.hanser-fachbuch.de/update)





Dani Schnider  
Claus Jordan  
Peter Welker  
Joachim Wehner

# **Data Warehouse Blueprints**

Business Intelligence  
in der Praxis

HANSER

Alle in diesem Buch enthaltenen Informationen, Verfahren und Darstellungen wurden nach bestem Wissen zusammengestellt und mit Sorgfalt getestet. Dennoch sind Fehler nicht ganz auszuschließen. Aus diesem Grund sind die im vorliegenden Buch enthaltenen Informationen mit keiner Verpflichtung oder Garantie irgendeiner Art verbunden. Autoren und Verlag übernehmen infolgedessen keine juristische Verantwortung und werden keine daraus folgende oder sonstige Haftung übernehmen, die auf irgendeine Art aus der Benutzung dieser Informationen – oder Teilen davon – entsteht.

Ebenso übernehmen Autoren und Verlag keine Gewähr dafür, dass beschriebene Verfahren usw. frei von Schutzrechten Dritter sind. Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Buch berechtigt deshalb auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.



Bibliografische Information der Deutschen Nationalbibliothek

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über [<http://dnb.d-nb.de>](http://dnb.d-nb.de) abrufbar.

Dieses Werk ist urheberrechtlich geschützt.

Alle Rechte, auch die der Übersetzung, des Nachdrucks und der Vervielfältigung des Buches, oder Teilen daraus, sind vorbehalten. Kein Teil des Werkes darf ohne schriftliche Genehmigung des Verlages in irgendeiner Form (Fotokopie, Mikrofilm oder ein anderes Verfahren), auch nicht für Zwecke der Unterrichtsgestaltung, reproduziert oder unter Verwendung elektronischer Systeme verarbeitet, vervielfältigt oder verbreitet werden.

© 2016 Carl Hanser Verlag München, [www.hanser-fachbuch.de](http://www.hanser-fachbuch.de)

Lektorat: Sylvia Hasselbach

Herstellung: Irene Weilhart

Copy editing: Sandra Gottmann, Münster-Nienberge

Umschlagdesign: Marc Müller-Bremer, [www.rebranding.de](http://www.rebranding.de), München

Umschlagrealisation: Stephan Rönigk

Gesamtherstellung: Kösel, Krugzell

Ausstattung patentrechtlich geschützt. Kösel FD 351, Patent-Nr. 0748702

Printed in Germany

Print-ISBN: 978-3-446-45075-2

E-Book-ISBN: 978-3-446-45111-7

# Inhalt

<b>Geleitwort</b> .....	<b>XIII</b>
<b>Über dieses Buch</b> .....	<b>XV</b>
<b>Die Autoren</b> .....	<b>XVII</b>
<b>1 Einleitung</b> .....	<b>1</b>
1.1 Ziele dieses Buches .....	2
1.2 Struktur dieses Buches .....	3
1.3 Hinweis zur Anwendung dieses Buches .....	4
<b>2 Architektur</b> .....	<b>5</b>
2.1 Data Warehouse-Architektur .....	5
2.1.1 Aufbau eines Data Warehouse .....	6
2.1.2 Transformationsschritte .....	9
2.1.3 Architekturgrundsätze .....	10
2.2 Architektur BI-Anwendungen .....	13
2.2.1 Die BI-Plattform zur Integration von Datenquellen .....	15
2.2.2 Die BI-Plattform zur Vereinheitlichung der Frontends .....	17
2.3 Datenhaltung .....	18
2.3.1 Grenzen gängiger DWH/BI-Technologien .....	19
2.3.2 Datenhaltung im Hadoop-Ecosystem .....	20
2.3.3 In-Memory-Datenbanken .....	23
<b>3 Datenmodellierung</b> .....	<b>27</b>
3.1 Vorgehensweise .....	27
3.1.1 Anforderungsgetriebene Modellierung .....	27
3.1.2 Quellsystemgetriebene Modellierung .....	29
3.1.3 Kombination der Ansätze .....	30
3.2 Relationale Modellierung .....	30
3.2.1 Darstellung von relationalen Datenmodellen .....	31
3.2.2 Normalisierung .....	31

3.2.3	Stammdaten und Bewegungsdaten .....	32
3.2.4	Historisierung .....	32
3.2.5	Relationales Core .....	34
3.2.6	Corporate Information Factory .....	35
3.2.7	Data Vault Modeling .....	35
3.3	Dimensionale Modellierung .....	37
3.3.1	Implementierung von dimensionalen Modellen .....	38
3.3.1.1	Relationale Implementierung .....	39
3.3.1.2	Multidimensionale Implementierung .....	40
3.3.2	Dimensionen .....	41
3.3.2.1	Fachliche Attribute .....	41
3.3.2.2	Technische Attribute .....	41
3.3.2.3	Hierarchien .....	42
3.3.2.4	Conformed Dimensions .....	43
3.3.2.5	Slowly Changing Dimensions .....	44
3.3.2.6	Zeitdimension .....	47
3.3.2.7	Bridge Tables .....	48
3.3.2.8	Spezielle Dimensionen .....	50
3.3.3	Fakten .....	51
3.3.3.1	Kennzahlen .....	51
3.3.3.2	Typen von Fakten .....	51
3.3.4	Modellierung spezieller Problemstellungen .....	53
3.3.4.1	Fakten unterschiedlicher Granularität und Rollen .....	53
3.3.4.2	Gemeinsame Hierarchiestufen in verschiedenen Dimensionen .....	54
3.3.4.3	Modellierungsgrundsätze für Dimensionen und Fakten .....	55
3.3.5	Darstellung von dimensionalen Modellen .....	56
3.3.5.1	ADAPT-Notation .....	56
3.3.5.2	Entity-Relationship-Diagramme .....	58
3.3.5.3	Data-Warehouse-Bus-Matrix .....	58
3.3.6	Dimensionales Core .....	59
3.4	Tools zur Datenmodellierung .....	60
3.4.1	Tools für relationale Datenmodellierung .....	60
3.4.2	Tools für dimensionale Datenmodellierung .....	61
<b>4</b>	<b>Datenintegration .....</b>	<b>63</b>
4.1	Data Profiling .....	64
4.1.1	Probleme mangelnder Datenqualität .....	64
4.1.2	Einsatz von Data Profiling .....	65
4.2	ETL .....	66
4.2.1	Aufgaben der ETL-Prozesse .....	67
4.2.1.1	Extraktion aus Quellsystemen .....	67
4.2.1.2	Transformationen .....	67
4.2.1.3	Laden in die Zieltabellen .....	68

4.2.2	ETL-Tools	68
4.2.2.1	Funktionalität von ETL-Tools	70
4.2.2.2	ETL oder ELT?	70
4.2.2.3	Positionierung von ETL-Tools	72
4.2.3	Performance-Aspekte	72
4.2.3.1	Mengenbasierte statt datensatzbasierte Verarbeitung	72
4.2.3.2	ELT-Tool statt ETL-Tool	73
4.2.3.3	Reduktion der Komplexität	74
4.2.3.4	Frühzeitige Mengeneinschränkung	75
4.2.3.5	Parallelisierung	76
4.2.4	Steuerung der ETL-Prozesse	78
4.2.4.1	Protokollierung des ETL-Ablaufs	78
4.2.4.2	Restartfähigkeit und Wiederaufsetzpunkte	79
4.3	Extraktion und Delta-Ermittlung	80
4.3.1	Delta-Extraktion im Quellsystem	81
4.3.1.1	Änderungsmarker und Journaltabellen	81
4.3.1.2	Delta-Ermittlung und Pending Commits	82
4.3.1.3	Change Data Capture	83
4.3.2	Voll-Extraktion und Delta-Abgleich im Data Warehouse	84
4.3.2.1	Zwei Versionen des Vollabzugs in der Staging Area	85
4.3.2.2	Vorteil einer Voll-Extraktion für die Delta-Ermittlung	87
4.3.3	Wann verwende ich was?	87
4.4	Fehlerbehandlung	88
4.4.1	Fehlende Attribute	89
4.4.1.1	Filtern von fehlerhaften Datensätzen	89
4.4.1.2	Fehlerhafte Datensätze in Fehlertabelle schreiben	89
4.4.1.3	Singletons auf Attributebene	90
4.4.2	Unbekannte Codewerte	90
4.4.2.1	Filtern von fehlerhaften Datensätzen	91
4.4.2.2	Singletons auf Datensatzebene	91
4.4.2.3	Generierung von Embryo-Einträgen	91
4.4.3	Fehlende Dimensionseinträge	92
4.4.3.1	Filtern von unvollständigen Fakten	93
4.4.3.2	Referenz auf Singleton-Einträge	94
4.4.3.3	Generieren von Embryo-Einträgen	95
4.4.4	Doppelte Datensätze	96
4.4.4.1	Verwendung von DISTINCT	97
4.4.4.2	Nur ersten Datensatz übernehmen	97
4.5	Qualitätschecks	97
4.5.1	Qualitätschecks vor und während des Ladens	98
4.5.2	Qualitätschecks nach dem Laden	99
4.5.3	Qualitätschecks mithilfe von Test-Tools	99
4.6	Real-Time BI	100
4.6.1	Begriffsbestimmung	101



4.6.2	Garantierte Verfügbarkeit von Informationen zu gegebenem Zeitpunkt	101
4.6.3	Verfügbarkeit von Informationen simultan zur Entstehung	102
4.6.4	Verfügbarkeit von Informationen kurz nach ihrer Entstehung	104
4.6.4.1	Events und Batchverarbeitung	105
4.6.4.2	Real-Time-Partitionen	106
4.6.5	Zusammenfassung	107
<b>5</b>	<b>Design der DWH-Schichten</b>	<b>109</b>
5.1	Staging Area	110
5.1.1	Gründe für eine Staging Area	111
5.1.2	Struktur der Stage-Tabellen	112
5.1.3	ETL-Logik für Stage-Tabellen	113
5.1.3.1	Einschränkungen bei der Extraktion	114
5.1.3.2	Transformation	114
5.1.3.3	Sonstige Informationen	115
5.2	Cleansing Area	115
5.2.1	Gründe für eine Cleansing Area	115
5.2.2	Struktur der Cleanse-Tabellen	116
5.2.3	Beziehungen in der Cleansing Area	118
5.2.4	ETL-Logik für Cleanse-Tabellen	120
5.2.4.1	Einschränkungen bei der Extraktion	121
5.2.4.2	Transformation	121
5.2.4.3	Sonstige Informationen	122
5.3	Core-Datenmodell allgemein	122
5.3.1	Aufgaben und Anforderungen an das Core	123
5.3.2	Stammdaten im Core	124
5.3.3	Bewegungsdaten im Core	124
5.3.4	Beziehungen im Core	124
5.3.5	Datenmodellierungsmethoden für das Core	125
5.4	Core-Datenmodell relational mit Kopf- und Versionstabellen	126
5.4.1	Historisierung von Stammdaten mit Kopf- und Versionstabellen	127
5.4.2	Struktur der Stammdatentabellen	128
5.4.2.1	Tabellenspalten und Schlüssel	129
5.4.2.2	Beziehungen (1:n) zwischen Stammdaten	132
5.4.2.3	Beziehungen (m:n) zwischen Stammdaten	133
5.4.3	ETL-Logik für Stammdatentabellen	135
5.4.3.1	Lookups (Schritt 1)	136
5.4.3.2	Outer Join (Schritt 2)	137
5.4.3.3	Neue Datensätze (Schritt 3)	141
5.4.3.4	Schließen einer Version/Fall 1 (Schritt 4)	142
5.4.3.5	Aktualisieren/Fall 2 (Schritt 5)	142
5.4.3.6	Versionieren/Fall 3 und 4 (Schritt 6)	142
5.4.3.7	Singletons	142

5.4.4	Typen von Bewegungsdaten .....	143
5.4.4.1	Transaction Tables .....	144
5.4.4.2	Snapshot Tables .....	144
5.4.4.3	Snapshot Tables versioniert .....	145
5.4.5	Struktur der Bewegungstabellen .....	146
5.4.5.1	Tabellenspalten und Schlüssel .....	147
5.4.5.2	Beziehungen zu Stammdaten .....	150
5.4.6	ETL-Logik für Bewegungstabellen .....	153
5.4.6.1	Lookups .....	154
5.4.6.2	Sonstige Informationen .....	155
5.4.7	Views für externen Core-Zugriff .....	155
5.4.7.1	Views für Stammdaten .....	156
5.4.7.2	Views für Bewegungsdaten .....	160
5.5	Core-Datenmodell relational mit Data Vault .....	161
5.5.1	Stammdaten .....	161
5.5.2	Beziehungen .....	162
5.5.3	Bewegungsdaten .....	162
5.5.4	Historisierung .....	163
5.5.5	Struktur der Tabellen .....	163
5.5.5.1	Hubtabellen – Tabellenspalten und Schlüssel .....	163
5.5.5.2	Satellitentabellen – Tabellenspalten und Schlüssel .....	164
5.5.5.3	Linktabellen – Tabellenspalten und Schlüssel .....	165
5.5.6	ETL-Logik .....	166
5.5.7	Views für externen Core-Zugriff auf das Data-Vault-Datenmodell .....	167
5.5.7.1	Views für Stammdaten (ein Satellite pro Hub bzw. Link) .....	167
5.5.7.2	Views für Stammdaten (mehrere Satellites pro Hub bzw. Link) .....	170
5.6	Core-Datenmodell dimensional .....	173
5.6.1	Star- oder Snowflake-Schema .....	174
5.6.1.1	Star-Schema .....	174
5.6.1.2	Snowflake-Schema .....	175
5.6.2	Historisierung von Stammdaten mit SCD .....	177
5.6.3	Struktur der Dimensionstabellen (Snowflake) .....	180
5.6.3.1	Tabellenspalten und Schlüssel .....	181
5.6.3.2	Beziehungen zwischen Hierarchiestufen .....	184
5.6.4	ETL-Logik für Dimensionstabellen (Snowflake) .....	185
5.6.4.1	Lookup .....	185
5.6.4.2	Weitere Schritte .....	186
5.6.5	Struktur der Faktentabellen (Snowflake) .....	186
5.6.6	ETL-Logik für Faktentabellen (Snowflake) .....	188
5.6.7	n:m-Beziehungen im dimensionalen Core .....	188
5.7	Marts .....	190
5.7.1	ROLAP oder MOLAP? .....	191
5.7.2	Historisierung von Data Marts .....	192
5.7.3	Star- oder Snowflake-Schema (ROLAP) .....	193

5.7.4	Struktur der Dimensionstabellen (Star)	194
5.7.4.1	Tabellenspalten und Schlüssel	194
5.7.4.2	Beispiel für Conformed Rollup	197
5.7.4.3	Beispiel für Dimension mit mehreren Hierarchien	198
5.7.5	ETL-Logik für Dimensionstabellen (Star)	199
5.7.5.1	Extraktion aus dem relationalen Core	200
5.7.5.2	Extraktion aus dem dimensional Core	207
5.7.6	Struktur der Faktentabellen (Star-Schema)	209
5.7.7	ETL-Logik für Faktentabellen (Star)	210
5.7.8	Multidimensionale Data Marts	210
5.7.8.1	Dimensionen (Cube)	211
5.7.8.2	Fakten (Cube)	212
<b>6</b>	<b>Physisches Datenbankdesign</b>	<b>215</b>
6.1	Indexierung	216
6.1.1	Staging Area	217
6.1.2	Cleansing Area	217
6.1.3	Core	217
6.1.4	Data Marts	218
6.2	Constraints	219
6.2.1	Primary Key Constraints	219
6.2.2	Foreign Key Constraints	220
6.2.3	Unique Constraints	221
6.2.4	Check Constraints	221
6.2.5	NOT NULL Constraints	222
6.3	Partitionierung	222
6.3.1	Grundprinzip von Partitionierung	223
6.3.2	Gründe für Partitionierung	223
6.3.3	Partitionierung in Staging und Cleansing Area	224
6.3.4	Partitionierung im Core	225
6.3.5	Partitionierung in den Data Marts	225
6.4	Datenkomprimierung	226
6.4.1	Redundanz	227
6.4.2	Wörterbuchmethode/Tokenbasierte Reduktion	227
6.4.3	Entropiekodierung	227
6.4.4	Deduplikation	228
6.4.5	Komprimierung bei spaltenorientierter Datenhaltung	228
6.5	Aggregationen	229
6.5.1	Vorberechnete Aggregationen	230
6.5.2	Query Rewrite	230
6.5.3	Einsatz im Data Warehouse	231

<b>7</b>	<b>BI-Anwendungen</b>	<b>233</b>
7.1	Überblick	233
7.2	Standardberichte	236
7.3	Ad-hoc-Analyse	238
7.4	BI-Portale	239
<b>8</b>	<b>Betrieb</b>	<b>241</b>
8.1	Release-Management	241
8.1.1	Kategorisierung der Anforderungen	242
8.1.2	Schnittstellen zu Quellsystemen	243
8.1.3	Umgang mit historischen Daten	245
8.1.4	Datenbankumgebungen	246
8.2	Deployment	248
8.2.1	Manuelles Deployment	248
8.2.2	Filebasiertes Deployment	249
8.2.3	Repository-basiertes Deployment	250
8.2.4	Kombiniertes Deployment	250
8.3	Monitoring	252
8.3.1	Betriebsmonitoring	252
8.3.2	System und DB-Monitoring	252
8.3.3	ETL-Monitoring	252
8.3.4	Performance-Monitoring	253
8.4	Migration	255
8.4.1	Datenbank	256
8.4.2	ETL-Tool	257
8.4.3	BI-Tools	258
	<b>Literatur</b>	<b>259</b>
	<b>Index</b>	<b>261</b>



# Geleitwort

*Von Dr. Carsten Bange, Gründer und Geschäftsführer des Business Application Research Centers (BARC), Teil des europäischen Analystenhauses CXP Group.*

Noch ein Buch über Data Warehousing? Ist darüber in den vergangenen 25 Jahren nicht genug geschrieben worden? Ich gebe zu, ich war skeptisch als die Autoren mich baten, ein Vorwort zu verfassen. Insbesondere auch, da wir in unserer täglichen Praxis als Marktanalysten eine deutlich wachsende Kritik vieler Unternehmen an ihrem Data Warehouse wahrnehmen. Insbesondere die Anwender verlangen nach Änderungen, um ihren veränderten Anforderungen Rechnung zu tragen. Die letzte BARC-Anwenderbefragung zu diesem Thema<sup>1</sup> zeigt deutlich, was den Veränderungsbedarf treibt: 62% der 323 befragten BI- und Data-Warehouse-Verantwortlichen sehen sich mit deutlich erhöhten Erwartungen in den Fachbereichen konfrontiert, 51% verstehen dabei die schnellere Veränderung von Geschäftsprozessen als wesentlichen Treiber für Anpassungen an Datenmanagement-Konzepten und 45% erfahren eine Unzufriedenheit mit der benötigten Zeit, um neue Anforderungen im Data Warehouse umzusetzen.

Vielen Unternehmen wird also immer klarer, dass sie die etablierten Data-Warehouse-Systeme so nicht mehr weiterbetreiben können, sondern hinsichtlich der Prozesse und der Organisation, IT-Architektur und eingesetzte Technologien und Werkzeuge komplett überdenken müssen.

Das vorliegende Buch liefert hierzu einen guten Beitrag und legt seinen Fokus dabei auf Methodik und Technologie. Es trägt eine große Menge von Erfahrungen zu „Best Practice“-Anleitungen zusammen, die helfen, das eigene Projekt auf eine solide Basis zu stellen und typische Fehler zu vermeiden. Es behandelt dabei auch neue Technologien, z.B. aus dem Hadoop- und NoSQL-Umfeld, die eine interessante Ergänzung der etablierten und ausgereiften Datenbank- und Datenintegrationstechnologien sein können. Die Autoren bieten damit Entwicklern, BI- und Data-Warehouse-Verantwortlichen ein solides methodisches Fundament, um die Informationsversorgung zur Entscheidungsfindung in Unternehmen erfolgreich aufzubauen.

Es bleibt dann im Unternehmen die wichtige Aufgabe, die verfügbaren Technologien in eine anforderungsgerechte Organisation einzubetten. Gerade Agilität und Flexibilität sind hier

---

<sup>1</sup> s. BARC-Anwenderbefragung „Modernes Datenmanagement für die Analytik“ (BARC 2015), Ergebnisstudie kostenfrei verfügbar unter [www.barc.de](http://www.barc.de) im Bereich Research.

die wesentlichen Anforderungen, die in den letzten Jahren beispielsweise den Trend zu „Self Service BI“ angefeuert haben, also der Bereitstellung weitgehender Möglichkeiten zur Zusammenstellung, Aufbereitung und Visualisierung von Daten für Fachanwender. Da dies häufig auch mit einer „Self Service-Datenintegration“ verbunden ist, ergibt sich schnell die Kehrseite der Medaille solcher Initiativen: Die Konsistenz von Daten kann in einer dezentralisierten Welt individueller Datenaufbereitung – wenn überhaupt – nur mit erheblichen Anstrengungen einer Data Governance sichergestellt werden. Die ersten Unternehmen kehren demnach auch schon wieder zu stärker zentralistisch ausgerichteten Konzepten zurück, um dem Daten-Wildwuchs Einhalt zu gebieten.

Dieser Spagat zwischen der Bereitstellung qualitätsgesicherter Daten unter übergreifender Kontrolle auf der einen sowie Flexibilität und Individualität in Datenzusammenstellung und -auswertung auf der anderen Seite ist aus unserer Sicht die momentan größte Herausforderung für Betreiber entscheidungsunterstützender Informationssysteme.

Das Buch zeigt, dass viele Methoden und Technologien hierfür zur Verfügung stehen. Werden sie richtig eingesetzt, sind dem Data Warehouse auch weitere 25 Jahre erfolgreichen Einsatzes beschieden, denn Entscheidungsträger im Unternehmen werden auch in Zukunft nicht auf konsistente und qualitätsgesicherte Daten zur Entscheidungsfindung verzichten.

Würzburg, den 14.7.2016

*Dr. Carsten Bange*

# Über dieses Buch

Das vorliegende Buch ist eine Weiterentwicklung des Buches „Data Warehousing mit Oracle – Business Intelligence in der Praxis“, das 2011 beim Carl Hanser Verlag erschienen und mittlerweile vergriffen ist. Im Vergleich zur vorherigen Version wurden hier die allgemeinen Konzepte, Architekturvorschläge und Vorgehensweisen stark ausgebaut und aktualisiert. Oracle-spezifische Informationen wurden – bis auf die Verwendung in Beispielen – weitgehend verallgemeinert, sodass die vorliegenden Blueprints auch für andere Datenbanktechnologien eingesetzt werden können.

Die Data Warehouse Blueprints wurden vorerst als interner Leitfaden für die BI-Consultants bei Trivadis zur Verfügung gestellt, bevor sie öffentlich publiziert wurden. Während dieser Zeit haben verschiedene Trivadis-Kollegen die einzelnen Kapitel überprüft und zahlreiche Korrekturen, Änderungsvorschläge und Ergänzungen zur nun vorliegenden Ausgabe beigetragen.





# Die Autoren

## Dani Schnider

Dani Schnider ist seit seinem abgeschlossenen Informatikstudium an der ETH Zürich (1990) in der Informatik tätig. Seit 1997 arbeitet er vorwiegend in DWH-Projekten. Konzeption, Design, Aufbau und Weiterentwicklung von Data Warehouses, logische und physische Datenmodellierung im DWH-Umfeld sowie Reviews, Architekturberatungen und Schulungen bilden seine Aufgabenschwerpunkte in diesem Bereich. Präsentationen an verschiedenen Konferenzen und Publikationen von Fachartikeln und Blog-Posts runden seine Tätigkeiten ab. (Kontakt: [dani.schnider@trivadis.com](mailto:dani.schnider@trivadis.com))

## Claus Jordan

Seit seinem Abschluss des Studiums der Wirtschaftsinformatik 1993 ist Claus Jordan im Umfeld Data Warehouse und Business Intelligence aktiv. Seit 2003 bringt er seine Erfahrung in diesen Bereichen für die Trivadis GmbH in zahlreichen Kundenprojekten, als Trainer und als Autor ein. Seine Schwerpunkte liegen dabei im Design unterschiedlicher Datenmodellierungsmethoden eines Data Warehouse, sowie in der Implementierung von ETL-Prozessen und deren Standardisierung. (Kontakt: [claus.jordan@trivadis.com](mailto:claus.jordan@trivadis.com))

## Peter Welker

Peter Welker arbeitete bereits vor dem Abschluss seines Studiums der Medizininformatik 1996 als Entwickler für Anwendungssoftware. 1998 wechselte er ins Data Warehousing und ist seitdem hauptsächlich in Projekten mit dem Fokus ETL, DWH-Lösungsarchitektur, Review und Performance aktiv. In den letzten Jahren beschäftigt er sich intensiv mit den neuen Technologien. Er präsentiert an Konferenzen, publiziert Fachartikel und verantwortet bei der Deutschen Oracle-Anwendergruppe (DOAG) das Thema „Big Data“. (Kontakt: [peter.welker@trivadis.com](mailto:peter.welker@trivadis.com))

## Joachim Wehner

Seit seiner Diplomarbeit „Werkzeuge zum Aufbau eines Data Warehouses“ aus dem Jahre 1996 lässt ihn dieses Thema nicht mehr los. Als Berater und Trainer arbeitet Joachim Wehner über die Jahre primär in BI-/DWH-Kundenprojekten. Im Mittelpunkt stehen dabei fast immer die Architektur, das Design sowie Reviews solcher Data-Warehouse-Umgebungen. Inzwischen hat sich sein Verantwortungsbereich von der Technik auf die Managementseite verlagert. (Kontakt: [joachim.wehner@trivadis.com](mailto:joachim.wehner@trivadis.com))

## ■ Danksagung

Die Kapitel dieses Buches wurden von verschiedenen Trivadis-Consultants geprüft, korrigiert und mit wertvollen Ergänzungen und Änderungsvorschlägen angereichert. Der Dank für diese Reviewarbeit gilt folgenden Personen: Adrian Abegglen, Aron Hennerdal, Beat Flühmann, Christoph Hisserich, Kamilla Reichardt, Maurice Müller, Peter Denk, Stanislav Lando, Thomas Brunner und Willfried Färber. Die gute und konstruktive Zusammenarbeit mit Frau Hasselbach und Frau Weilhart vom Hanser Verlag ist an dieser Stelle ebenfalls dankend zu erwähnen wie das passende Geleitwort zu diesem Buch von Herrn Dr. Carsten Bange vom Business Application Research Center (BARC).

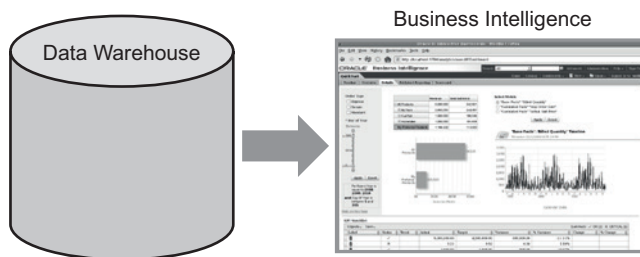
Ein besonderer Dank gilt natürlich der Trivadis AG für die Unterstützung des Buchprojekts während der Erstellung und dafür, dass sie es ermöglicht hat, dass das Buch wiederum öffentlich publiziert werden kann.

# 1

## Einleitung

Business Intelligence (BI) und Data Warehousing (DWH) sind zwei Begriffe, die in vielen Unternehmen nicht mehr wegzudenken sind und denen eine immer wichtigere Bedeutung zukommt. In vielen großen Unternehmen gehören Data Warehouses und BI-Applikationen zu den zentralen Systemen. Auch kleinere und mittlere Betriebe benutzen Business Intelligence für die Planung und Überprüfung ihrer Geschäftsziele.

Business Intelligence bezeichnet die systematische Auswertung von Daten eines Unternehmens, um damit Geschäftsprozesse zu analysieren und zu optimieren. Das Ziel von Business Intelligence ist es, aus vergleichbaren Kennzahlen neue Erkenntnisse zu gewinnen. Sie dienen als Basis für strategische und operative Entscheidungen, mit denen die Unternehmensziele besser erreicht werden können.



**Bild 1.1** Data Warehouse als Basis für Business Intelligence

Um Business Intelligence erfolgreich betreiben zu können, muss eine solide Data-Warehouse-Architektur als Basis vorhanden sein. Ziel eines Data Warehouse ist es, die Daten aus verschiedenen operativen Quellsystemen so zusammenzuführen und in geeigneter Form abzuspeichern, dass darauf einfache und flexible Abfragen sowie verschiedenste Arten von Auswertungen möglich sind.



### Data Warehouse oder Business Intelligence?

Die Begriffe „Data Warehouse“ und „Business Intelligence“ werden teilweise in unterschiedlichem Kontext verwendet. So wird „Business Intelligence“ einerseits als Sammelbegriff für BI-Gesamtlösungen – unter anderem Data Warehouses – verwendet, andererseits für die systematische Analyse von Geschäftsprozessen anhand von Kennzahlen.

Ein Data Warehouse wiederum bezeichnet zum einen ein System aus Datenbank(en) und ETL-Prozessen, das die notwendigen Informationen zur Verfügung stellt, um Business Intelligence betreiben zu können. Zum anderen wird der Begriff „Data Warehouse“ oft auch für die zentrale Integrations- und Historisierungsschicht innerhalb eines DWH-Gesamtsystems verwendet.

Eine einheitliche Namensgebung innerhalb der Informatikwelt ist kaum möglich, aber zumindest im vorliegenden Buch wurde versucht, die Begriffe einheitlich zu verwenden. Hier verstehen wir unter Business Intelligence (BI) jede Art von Anwendungen zur Datenanalyse, basierend auf einem Data Warehouse (DWH). Das Data Warehouse ist somit die technische Datenbasis für Business Intelligence. Die zentrale Schicht im DWH wird hier – zur Unterscheidung vom DWH-Gesamtsystem – als Core (oder Core Data Warehouse) bezeichnet.

Ein Data Warehouse umfasst somit die technische und fachliche Basis, die notwendig ist, um Anwendungen im Bereich Business Intelligence betreiben zu können. Integration, Skalierbarkeit und Performance sind wichtige Erfolgsfaktoren. Jede BI-Applikation und jedes Data Warehouse kann nur dann erfolgreich sein, wenn die Architektur richtig aufgebaut ist, die einzelnen Komponenten zusammenpassen und das Gesamtsystem fehlerlos konfiguriert wird.

In den Data Warehouse Blueprints werden Grundlagen und Konzepte für den Aufbau und Betrieb von Data Warehouses beschrieben. Die vorliegenden Kapitel wurden soweit möglich unabhängig von einer spezifischen Technologie beschrieben und lassen sich mit unterschiedlichen Datenbanksystemen und Softwarekomponenten umsetzen. Da jedoch die Autoren alle im Oracle-Umfeld tätig sind, schimmern teilweise technologiespezifische Detailinformationen durch. Die mit anderen Datenbanktechnologien vertrauten Leserinnen und Leser werden aufgefordert, beim Lesen der folgenden Kapitel tolerant zu sein und die technologiespezifischen Begriffe sinngemäß in ihre Nomenklatur zu übersetzen.

## ■ 1.1 Ziele dieses Buches

Um leistungsfähige und stabile DWH-Systeme aufbauen und betreiben zu können, ist entsprechendes Know-how über Data Warehousing notwendig. Das vorliegende Buch gibt einen Überblick über eine typische DWH-Architektur und zeigt anhand von zahlreichen Beispielen auf, wie die einzelnen Komponenten eines Data Warehouse realisiert und betrie-

ben werden können. Der Hauptfokus liegt dabei nicht auf einer vollständigen Aufzählung der technischen Möglichkeiten – dazu stehen die Dokumentationen oder entsprechende Schulungen für die jeweilige Datenbanktechnologie zur Verfügung –, sondern darauf, wie allgemeine Technologien und Methoden in konkreten DWH-Projekten verwendet werden können. Die hier aufgezeigten Konzepte und Vorgehensweisen wurden in zahlreichen Projekten eingesetzt und – basierend auf den Erfahrungen daraus – verfeinert und erweitert.

Anhand verschiedener Tipps und Tricks aus der Praxis wird erläutert, wie die beschriebenen Methoden im Data Warehousing eingesetzt werden können. Es versteht sich von selbst, dass die hier beschriebenen Möglichkeiten keinen Anspruch auf Vollständigkeit erheben. Jedes Data Warehouse hat andere Anforderungen, Systemvorgaben und Spezialfälle, die zu berücksichtigen sind. Die hier vorgestellten Konzepte sollen jedoch als technischer Leitfaden dienen, um auch komplexe und umfangreiche Data Warehouses nach bewährtem Muster aufbauen zu können.

## ■ 1.2 Struktur dieses Buches

Das vorliegende Buch ist in folgende Hauptkapitel unterteilt:

- Kapitel 1 gibt einen Überblick über Data Warehousing und Business Intelligence, die Struktur des Buches und die verwendeten Begriffe.
- Kapitel 2 beschreibt die grundlegenden *Architekturen* von Data Warehouse, BI-Anwendungen und Datenhaltung innerhalb eines DWH-Systems.
- Kapitel 3 befasst sich mit der *Datenmodellierung* im Data Warehouse und beschreibt unterschiedliche Modellierungsansätze, wie sie im DWH-Umfeld zum Einsatz kommen.
- Kapitel 4 geht auf verschiedene Aspekte der *Datenintegration* von den Quellsystemen ins Data Warehouse und beschreibt verschiedene Konzepte, die in diesem Zusammenhang verwendet werden.
- Kapitel 5 befasst sich detailliert mit dem *Design der DWH-Schichten*, basierend auf der Architektur und den Grundsätzen der Datenmodellierung und Datenintegration der vorhergehenden Kapitel.
- Kapitel 6 beschreibt verschiedene Thematiken im Zusammenhang mit dem *physischen Datenbankdesign* einer DWH-Datenbank, ohne auf spezifische Features einzelner Datenbanksysteme einzugehen.
- Kapitel 7 gibt einen kurzen Überblick über verschiedene Kategorien von *BI-Anwendungen*, wie sie typischerweise in Business-Intelligence-Plattformen mit Data Warehouses zum Einsatz kommen.
- Kapitel 8 befasst sich mit unterschiedlichen Aspekten, die beim *Betrieb* eines Data Warehouse berücksichtigt werden müssen. Dazu gehören Themen wie Release Management, Deployment, Monitoring und Migration von Data Warehouses.

## ■ 1.3 Hinweis zur Anwendung dieses Buches

Als „Blueprints“ – also Baupläne – werden hier Verfahren und Methoden bezeichnet, die sich in verschiedenen DWH-Projekten bewährt haben. Das bedeutet aber nicht, dass sie die einzige oder beste Lösung für **jedes** Data Warehouse beschreiben.

Erfahrene DWH-Entwickler und BI-Consultants werden nach der Lektüre der nachfolgenden Kapitel in der Lage sein, die beschriebenen Konzepte und Praxistipps soweit sinnvoll in konkreten DWH-Projekten anzuwenden. Sie sollten aber auch die nötige Erfahrung haben, bei Bedarf zu erkennen, ob und wann von den hier beschriebenen Blueprints abzuweichen ist – sei es durch den Einsatz anderer Technologien, durch spezielle Kundenbedürfnisse oder durch eine andere Ausgangslage, als sie hier angenommen wird.

Das gilt auch für die Architektur eines Data Warehouse: Die im vorliegenden Buch verwendete Architektur beschreibt eine bewährte, aber nicht die einzig mögliche Variante, wie ein Data Warehouse auszusehen hat. Vielleicht hat das DWH bei einem spezifischen Kunden mehr oder weniger Schichten, oder es werden andere Begriffe für die einzelnen Komponenten verwendet. Auch hier gilt: Was zweckmäßig ist, kann und soll aus den Blueprints übernommen werden. Wenn es sinnvoll und begründbar ist, von den hier beschriebenen Methoden abzuweichen, ist dies durchaus erlaubt und gewünscht – sofern es der Qualität des zu bauenden Data Warehouse und der Zufriedenheit des Kunden dient.

# 2

## Architektur

Eine gut strukturierte Architektur ist eine wichtige Voraussetzung für den erfolgreichen Einsatz von Data Warehousing und Business Intelligence. Die wichtigsten Grundsätze zur Architektur von DWH-Systemen und BI-Anwendungen sowie verschiedene Möglichkeiten zur Datenhaltung in einem Data Warehouse sind in diesem Kapitel zusammengefasst.

- Abschnitt 2.1 beschreibt die grundlegende Architektur eines Data Warehouse und stellt die verschiedenen Schichten einer DWH-Architektur vor.
- In Abschnitt 2.2 werden die wichtigsten Grundsätze zur Architektur und zum Aufbau von BI-Anwendungen erläutert.
- Abschnitt 2.3 gibt einen Überblick über verschiedene Konzepte zur Datenhaltung, wie sie für Data Warehouses zum Einsatz kommen.

### ■ 2.1 Data Warehouse-Architektur

Ein Data Warehouse (DWH) stellt die technische Infrastruktur zur Verfügung, die benötigt wird, um Business Intelligence betreiben zu können. Sein Zweck ist es, Daten aus unterschiedlichen Datenquellen zu integrieren und eine historisierte Datenbasis zur Verfügung zu stellen, welche für Standard- und Ad-hoc-Reporting, OLAP<sup>1</sup>-Analysen, Balanced Scorecards, BI-Dashboards und weitere BI-Anwendungen eingesetzt werden kann. Ein DWH ist ein abfrageoptimiertes System, mit welchem auf eine Sammlung von historisierten Daten über einen längeren Zeitpunkt zugegriffen werden kann.

Durch diese Ausgangslage ergeben sich einige Unterschiede zwischen einem operativen System (auch OLTP<sup>2</sup>-System genannt) und einem Data Warehouse. Während in einem OLTP-System mehrere bis viele Anwender gleichzeitig Daten einfügen, ändern und löschen, ist dies bei einem DWH-System in der Regel nicht der Fall. Die einzigen „Anwender“, die in ein Data Warehouse schreiben, sind die ETL-Prozesse, welche Daten von den Quellsystemen ins DWH laden. Auch die Art der Abfragen ist unterschiedlich. In operativen Systemen werden typi-

---

<sup>1</sup> OLAP = Online Analytical Processing

<sup>2</sup> OLTP = Online Transaction Processing



scherweise spezifische Informationen in einem großen Datenbestand gesucht, beispielsweise die letzten zehn Banktransaktionen eines bestimmten Kunden. In einem Data Warehouse hingegen werden meistens Auswertungen über große Datenmengen ausgeführt und aggregiert, zum Beispiel die Summe über alle Verkäufe an alle Kunden einer bestimmten Region.

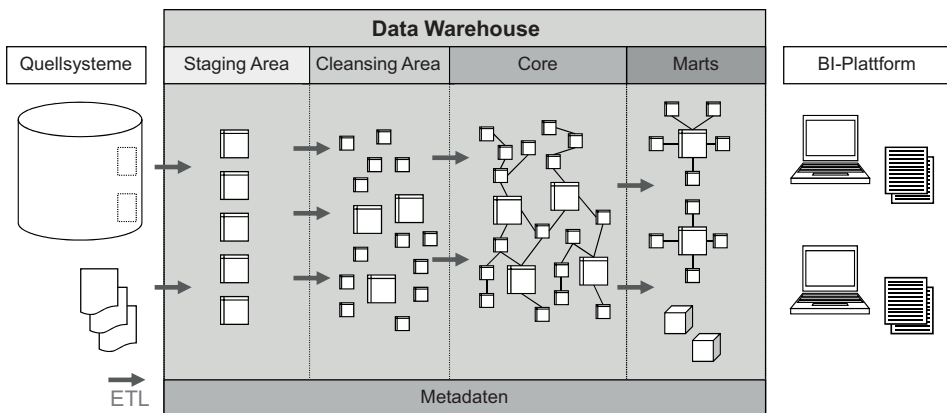
Um diesen unterschiedlichen Bedürfnissen gerecht zu werden, werden DWH-Datenbanken anders aufgebaut als OLTP-Datenbanken. Architektur, Design und Datenmodellierung funktionieren im DWH-Umfeld nicht auf die gleiche Weise, wie es viele erfahrene Architekten, Datenmodellierer und Entwickler gewohnt sind, die hauptsächlich im Bereich von OLTP-Datenbanken tätig sind. Auf die spezifischen Bedürfnisse von DWH-Systemen wird deshalb nachfolgend eingegangen.

Die Komplexität und Erweiterbarkeit eines Data Warehouse ist weitgehend abhängig von der verwendeten Architektur. Deshalb ist es in jedem DWH-Projekt von Anfang an wichtig, dass eine saubere Architektur definiert und implementiert wird. In der Regel bedeutet das, dass die Architektur aus unterschiedlichen Schichten besteht. Diese Schichten decken jeweils unterschiedliche Anforderungen ab. Auch wenn dies zum Beginn des Projektes nach Mehraufwand aussieht, zahlt sich eine konsequente Aufteilung in verschiedene DWH-Schichten im späteren Projektverlauf und im operativen Betrieb des Systems aus.

Leider wird oft der Fehler gemacht, dass aufgrund von knappen Terminvorgaben wesentliche Architekturgrundsätze missachtet und „Abkürzungen“ bzw. „Schnellschüsse“ implementiert werden. Diese Ad-hoc-Lösungen können früher oder später zu Problemen führen. Der Aufwand, diese wiederum zu beheben, ist oft größer als der Aufwand, von Anfang an eine saubere Lösung zu realisieren.

### 2.1.1 Aufbau eines Data Warehouse

Ein Data Warehouse besteht typischerweise aus verschiedenen Schichten (auch Layers, Bereiche oder Komponenten genannt) und Datenflüssen zwischen diesen Schichten. Auch wenn nicht jedes DWH-System alle Schichten umfassen muss, lässt sich jedes Data Warehouse auf eine Grundarchitektur, wie sie in Bild 2.1 dargestellt ist, zurückführen.



**Bild 2.1** Grundarchitektur eines Data Warehouse

*Um den Zweck der einzelnen Schichten in einer DWH-Architektur zu erklären, werden nachfolgend Beispiele aus dem „realen Leben“ gezeigt. Nehmen wir an, das DWH sei ein großes Lebensmittelgeschäft. Auch dort gibt es verschiedene Bereiche, die jeweils einem bestimmten Zweck dienen.*

Folgende Schichten oder Bereiche gehören zu einer vollständigen DWH-Architektur:

- **Staging Area:** Daten aus unterschiedlichen Quellsystemen werden zuerst in die Staging Area geladen. In diesem ersten Bereich des DWH werden die Daten so gespeichert, wie sie angeliefert werden. Die Struktur der Stage-Tabellen entspricht deshalb der Schnittstelle zum Quellsystem<sup>3</sup>. Beziehungen zwischen den einzelnen Tabellen bestehen keine. Jede Tabelle enthält die Daten der letzten Lieferung, welche vor der nächsten Lieferung gelöscht werden.

*In einem Lebensmittelgeschäft entspricht die Staging Area der Laderampe, an der die Lieferanten (Quellsysteme) ihre Waren (Daten) abliefern. Auch dort werden immer nur die neuesten Lieferungen zwischengelagert, bevor sie in den nächsten Bereich überführt werden.*

- **Cleansing<sup>4</sup> Area:** Bevor die gelieferten Daten ins Core geladen werden, müssen sie bereinigt werden. Fehlerhafte Daten müssen entweder ausgefiltert, korrigiert oder durch Singletons (Defaultwerte) ergänzt werden. Daten aus unterschiedlichen Quellsystemen müssen in eine vereinheitlichte Form transformiert und integriert werden. Die meisten dieser Bereinigungsschritte werden in der Cleansing Area durchgeführt. Auch diese Schicht enthält nur die Daten der letzten Lieferung.

*Im Lebensmittelgeschäft kann die Cleansing Area mit dem Bereich verglichen werden, in dem die Waren für den Verkauf kommissioniert werden. Die Waren werden ausgepackt, Gemüse und Salat werden gewaschen, das Fleisch portioniert, ggf. mehrere Produkte zusammengefasst und alles mit Preisetiketten versehen. Die Qualitätskontrolle der angelieferten Ware gehört ebenfalls in diesen Bereich.*

- **Core:** Die Daten aus den verschiedenen Quellsystemen werden über die Staging und Cleansing Area in einem zentralen Bereich, dem Core, zusammengeführt und dort über einen längeren Zeitraum, oft mehrere Jahre, gespeichert. Eine Hauptaufgabe des Core ist es, die Daten aus den unterschiedlichen Quellen zu integrieren und nicht mehr getrennt nach Herkunft, sondern themenspezifisch strukturiert zu speichern. Oft spricht man bei thematischen Teilbereichen im Core von „Subject Areas“. Die Daten werden im Core so abgelegt, dass historische Daten zu jedem späteren Zeitpunkt ermittelt werden können. Das Core sollte die einzige Datenquelle für die Data Marts sein. Direkte Zugriffe von Benutzern auf das Core sollten möglichst vermieden werden.

*Das Core kann mit einem Hochregallager verglichen werden. Waren werden so abgelegt, dass sie jederzeit auffindbar sind, aber der Zugriff darauf ist nur internen Mitarbeitern möglich. Kunden haben im Lager nichts zu suchen – außer vielleicht bei IKEA. Im Gegensatz zu einem Hochregallager bleiben die Daten aber auch dann im Core erhalten, nachdem sie an die Data Marts übertragen wurden.*

<sup>3</sup> Oft werden den Stage-Tabellen zusätzliche Attribute für Auditinformationen zugefügt, die im Quellsystem nicht vorhanden sind.

<sup>4</sup> Der Begriff „Cleansing“ wird englisch „klensing“ ausgesprochen und nicht „kliinsing“.

- **Marts:** In den Data Marts werden Teilmengen der Daten aus dem Core so aufbereitet abgespeichert, dass sie in einer für die Benutzerabfragen geeigneten Form zur Verfügung stehen. Jeder Data Mart sollte nur die für die jeweilige Anwendung relevanten Daten bzw. eine spezielle Sicht auf die Daten enthalten. Das bedeutet, dass typischerweise mehrere Data Marts für unterschiedliche Benutzergruppen und BI-Anwendungen definiert werden. Dadurch kann die Komplexität der Abfragen reduziert werden. Das erhöht die Akzeptanz des DWH-Systems bei den Benutzern.

*Die Data Marts sind die Marktstände oder Verkaufsgestelle im Lebensmittelgeschäft. Jeder Marktstand bietet eine bestimmte Auswahl von Waren an, z.B. Gemüse, Fleisch oder Käse. Die Waren werden so präsentiert, dass sie von der jeweiligen Kundengruppe akzeptiert, also gekauft werden.*

- **ETL-Prozesse:** Die Daten, die von den Quellsystemen als Files, Schnittstellentabellen oder über einen View Layer zur Verfügung gestellt werden, werden in die Staging Area geladen, in der Cleansing Area bereinigt und dann im Core integriert und historisiert. Vom Core werden aufgrund von fachlichen Anforderungen Teilmengen oder oft auch nur Aggregate in die verschiedenen Data Marts geladen.

All diese Datenflüsse werden unter dem Begriff ETL (Extraction, Transformation, Loading) zusammengefasst. Die Extraktion der Daten aus den Quellsystemen findet in der Regel außerhalb des DWH-Systems statt, nämlich in den Quellsystemen selbst. Als Transformationen werden alle Datenumformungen, Bereinigungen, Anreicherungen mit Zusatzinformationen und Aggregationen bezeichnet. Schließlich werden die Daten in die Zieltabellen der nächsten Schicht geladen.

*Die ETL-Prozesse sind die Mitarbeiter des Lebensmittelgeschäfts, die unterschiedliche Arbeiten verrichten müssen, damit die Lebensmittel vom Lieferanten bis hin zum Kunden gelangen.*

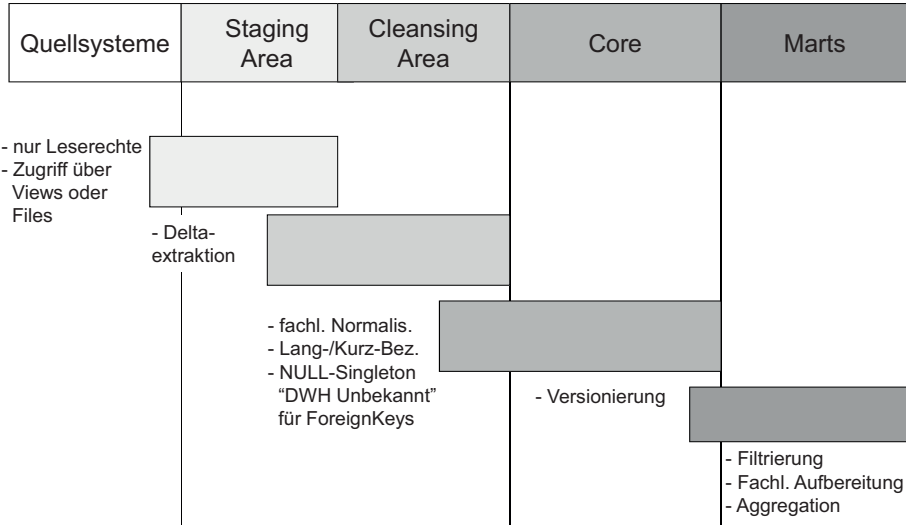
- **Metadaten:** Für den reibungsfreien Betrieb des Data Warehouse werden unterschiedliche Arten von Metadaten benötigt. Fachliche Metadaten enthalten fachliche Beschreibungen aller Attribute, Drill-Pfade und Aggregationsregeln für die Frontend-Applikationen und Codebezeichnungen. Technische Metadaten beschreiben z.B. Datenstrukturen, Mapping-Regeln und Parameter zur ETL-Steuerung. Operative Metadaten beinhalten alle Log-Tabellen, Fehlermeldungen, Protokollierungen der ETL-Prozesse und vieles mehr. Die Metadaten bilden die Infrastruktur eines DWH-Systems und werden als „Daten über Daten“ beschrieben.

*Auch in unserem Lebensmittelgeschäft braucht es eine funktionierende Infrastruktur – von Wegweisern zur Kasse bis hin zur Klimaüberwachung der Frischwaren.*

Nicht jedes Data Warehouse hat genau diesen Aufbau. Teilweise werden einzelne Bereiche zusammengefasst – zum Beispiel Staging Area und Cleansing Area – oder anders bezeichnet. So wird zum Teil das Core als „Integration Layer“ oder als „(Core) Data Warehouse“ bezeichnet. Wichtig ist jedoch, dass das Gesamtsystem in verschiedene Bereiche unterteilt wird, um die unterschiedlichen Aufgabenbereiche wie Datenbereinigung, Integration, Historisierung und Benutzerabfragen zu entkoppeln. Auf diese Weise kann die Komplexität der Transformationsschritte zwischen den einzelnen Schichten reduziert werden.

### 2.1.2 Transformationsschritte

Mithilfe der ETL-Prozesse werden die Daten von den Quellsystemen ins Data Warehouse geladen. In jeder DWH-Schicht werden dabei unterschiedliche Transformationsschritte durchgeführt, wie in Bild 2.2 dargestellt.



**Bild 2.2** Transformationsschritte im Data Warehouse

#### ■ *Quellsystem → Staging Area*

Wird direkt auf ein relationales Quellsystem zugegriffen, sollte als Schnittstelle zwischen Quellsystem und DWH ein View Layer als definierte Zugriffsschicht implementiert werden. Der Zugriff der ETL-Prozesse auf das Quellsystem erfolgt dann ausschließlich über diese Views. Auf diese Weise kann eine gewisse Unabhängigkeit der ETL-Prozesse gegenüber Strukturänderungen auf dem Quellsystem erreicht werden. Die Views können außerdem für die Delta-Extraktion verwendet werden, indem sie so implementiert werden, dass nur die jeweils relevante Teilmenge der Daten in den Views zur Verfügung steht. Dieses Verfahren wird teilweise für Change Data Capture (CDC) verwendet.

Als Alternative zum direkten Zugriff auf das Quellsystem werden häufig Dateien als Schnittstelle zwischen Quellsystem und Data Warehouse verwendet. Die Extraktion der Daten in die Dateien erfolgt auf dem Quellsystem und wird meistens außerhalb des DWH-Projektes realisiert.

Die Daten, ob über Views oder Files geliefert, werden unverändert in die Staging Area geladen und ggf. mit Auditinformationen angereichert.

#### ■ *Staging Area → Cleansing Area*

Beim Laden in die Cleansing Area werden die Daten geprüft, bereinigt und mit zusätzlichen Attributen angereichert. Dazu gehört zum Beispiel die Ermittlung von Lang- und Kurztexten aus den fachlichen Attributen der Quellsysteme. Fehlende oder fehlerhafte Attribute und Foreign Keys werden durch Singletons ersetzt.

Fehlerhafte Datensätze können je nach Anforderungen ignoriert, aufgrund von fixen Regeln korrigiert, durch Singletons ersetzt oder in Fehlertabellen geschrieben werden. Fehlertabellen können als Basis für Fehlerprotokolle oder manuelle Korrekturen verwendet werden. Bei solchen aufwendigen Varianten der Fehlerbehandlung muss allerdings organisatorisch geklärt werden, wer für die Fehlerkorrekturen verantwortlich ist.

#### ■ *Cleansing Area* → *Core*

Nachdem die Daten in der Cleansing Area in die benötigte Form aufbereitet wurden, werden sie ins Core geladen. In diesem Schritt findet die Versionierung der Stammdaten<sup>5</sup> statt, d. h., es wird für jeden Datensatz geprüft, ob sich etwas geändert hat und somit eine neue Version erstellt werden muss. Je nach Historisierungsanforderungen und Core-Datenmodell gibt es verschiedene Varianten der Versionierung von Stammdaten.

Die Bewegungsdaten<sup>6</sup> werden historisiert. Weil sich Bewegungsdaten nachträglich nicht mehr ändern, heißt das, dass laufend neue Daten eingefügt und über einen längeren Zeitraum gespeichert werden. Oft besteht die Anforderung, dass Bewegungsdaten nach einer gewissen Zeit – in der Regel nach mehreren Jahren – aus dem Core gelöscht werden.

Aggregationen werden im Core nicht durchgeführt. Die Bewegungsdaten werden auf der Detaillierungsstufe, die geliefert wird, gespeichert.

#### ■ *Core* → *Marts*

Die Transformationen vom Core in die Data Marts bestehen aus der Filtrierung der Daten auf die für jeden Data Mart erforderliche Teilmenge, der fachlichen Aufbereitung der Dimensionen<sup>7</sup> in die gewünschten Hierarchiestufen sowie – falls erforderlich – der Aggregation der Bewegungsdaten auf die Granularität der Faktentabellen.

### 2.1.3 Architekturgrundsätze

Obwohl sich die Architektur vieler DWH-Systeme in Details unterscheidet und oft auch unterschiedliche Namen für die einzelnen Bereiche verwendet werden, gibt es ein paar wichtige Architekturgrundsätze, die auf jeden Fall berücksichtigt werden sollten. Vereinfachungen der Architektur sind erlaubt, aber die wichtigsten Schichten sollten auf keinen Fall weggelassen werden.

<sup>5</sup> Stammdaten (oder Referenzdaten) sind zustandsorientierte Daten, die sich im Laufe der Zeit ändern können. Um die Änderungen im Core nachvollziehbar abspeichern zu können, werden die Daten versioniert. Das heißt, dass für jede Datenänderung ein neuer Datensatz im Core eingefügt wird. Die Versionierung von Stammdaten wird in Kapitel 3 (Datenmodellierung) genau erklärt.

<sup>6</sup> Bewegungsdaten (oder Transaktionsdaten) sind ereignisorientierte Daten, die aufgrund eines bestimmten Ereignisses (z. B. Transaktion, Messung) entstehen und nachträglich nicht mehr geändert werden. Sie sind immer mit einem Ereigniszeitpunkt (z. B. Transaktionsdatum) verbunden. Die Historisierung von Bewegungsdaten wird in Kapitel 3 (Datenmodellierung) beschrieben.

<sup>7</sup> Die Begriffe Dimensionen, Hierarchien und Fakten sind Elemente der dimensionalen Modellierung und werden in Kapitel 3 (Datenmodellierung) erläutert.