
Data Warehousing and Business Intelligence

Introduction

- Operational versus Tactical/Strategic Decision Making
- Data Warehouse Definition
- Data Warehouse Schemas
- The Extraction Transformation and Loading (ETL) Process
- Data Marts
- Virtual Data Warehouses and Virtual Data Marts
- Operational Data Store
- Data Warehouses versus Data Lakes
- Business intelligence

Decision Making Levels

- Operational level
 - Day-to-day business decisions
 - Short time frame
 - On-Line Transaction Processing (OLTP)
- Tactical level
 - Middle management decisions
 - Medium-term focus
- Strategic level
 - Senior management
 - Long-term focus

Decision Making Levels

- Operational systems
 - Operational level
 - Focus on simple INSERT, UPDATE, DELETE and/or SELECT statements
 - Transaction throughput
- Decision Support Systems
 - Tactical + strategic level
 - Focus on data retrieval by answering complex ad-hoc queries (SELECT statements)
 - Represent data in a multidimensional way
 - Trend analysis

Data Warehouse Definition

- *“A data warehouse is a subject-oriented, integrated, time-variant, and nonvolatile collection of data in support of management’s decision-making process.” (Inmon, 1996)*
- Subject oriented
 - organized around subjects such as customers, products, sales
- Integrated
 - integrates data from a variety of operational sources and a variety of formats

Data Warehouse Definition

- Non-volatile
 - data is primarily read-only
 - data loading and data retrieval
- Time variant
 - time series of periodic snapshots

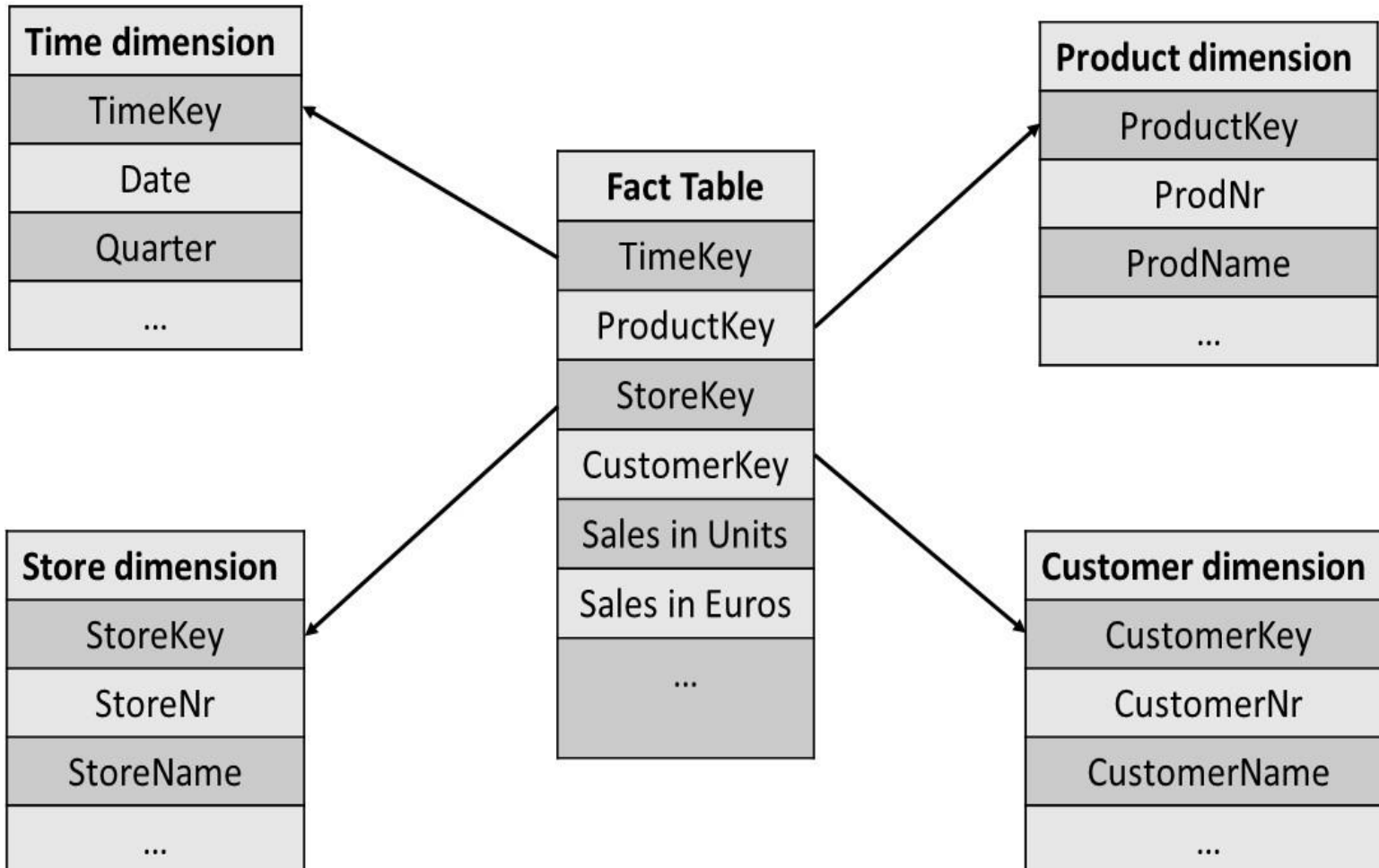
Data Warehouse Definition

	Transactional system	Data Warehouse
Usage	Day to day business operations	Decision support at tactical/strategic level
Data latency	Real-time data	Periodic snapshots, incl. historical data
Design	Application oriented	Subject Oriented
Normalization	Normalized data	(Sometimes also) denormalized data
Data manipulation	Insert/Update/Delete/Select	Insert/Select
Transaction management	Important	Less of a concern
Type of queries	Many, simple queries	Fewer, but complex and ad-hoc queries

Data Warehouse Schemas

- Star schema
- Snowflake schema
- Fact constellation
- Specific Schema issues

Star schema



Star schema

- Fact table
 - one tuple per transaction or event (i.e., a fact) and also measurement data (e.g. sales)
- Dimension table
 - further information about each of the facts in the fact table (e.g., Time, Store, Customer, Product).
 - criteria for aggregating the measurement data
 - often denormalized

Star Schema

Fact table

TimeKey	ProductKey	StoreKey	CustomerKey	Sales in Units	Sales in Euros
200	50	100	20006010	6	167.94
210	25	130	20006012	3	54
180	30	150	20006008	12	384
...					

Dimension tables

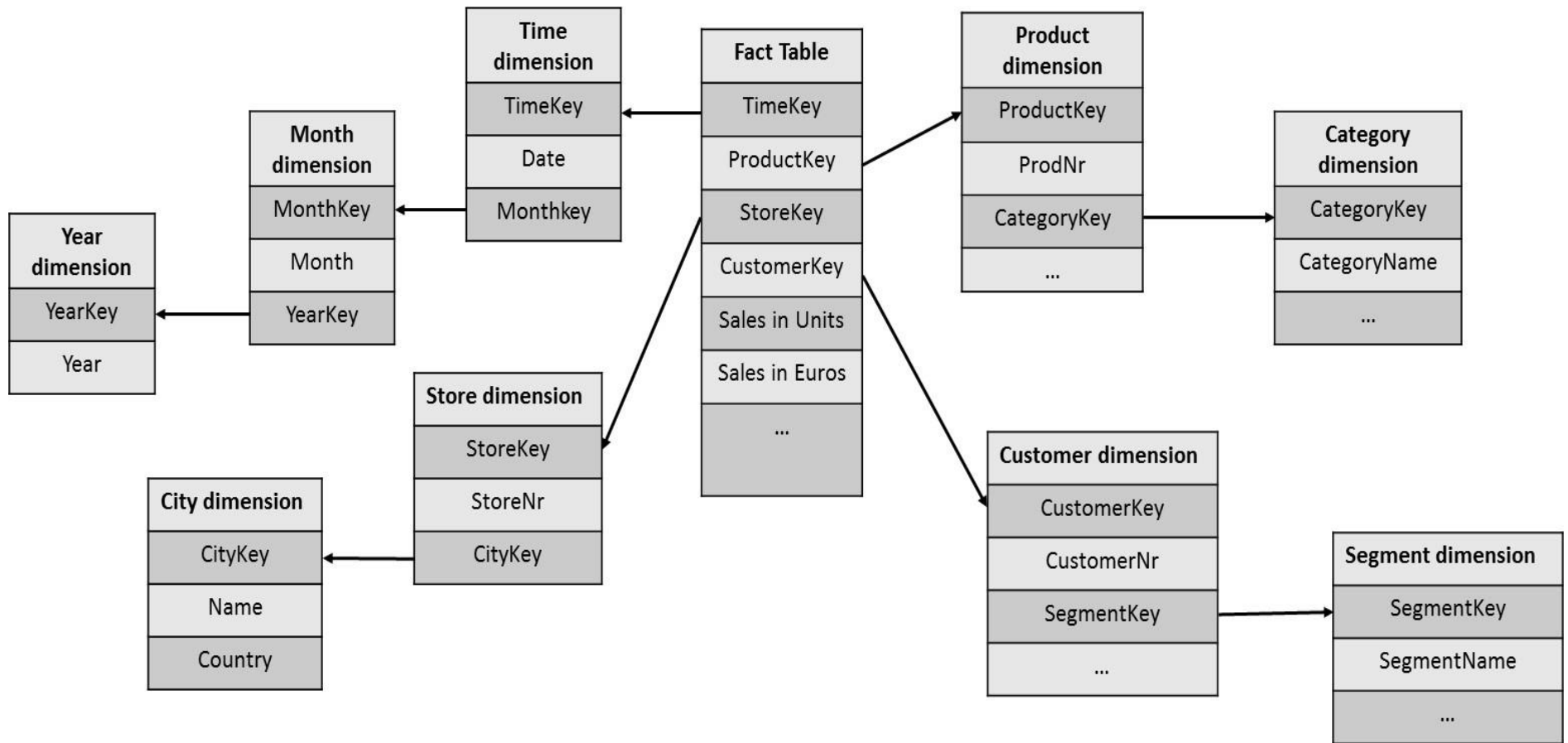
TimeKey	Date	Quarter	...
200	08/03/2017	1	...
210	09/11/2017	3	
180	27/02/2017	1	

CustomerKey	CustomerNr	CustName	...
20006008	20	Bart Baesens	...
20006010	10	Wilfried Lemahieu	
20006012	5	Seppe vanden Broucke	

StoreKey	StoreNr	StoreName	...
100	68	The Wine Depot	...
130	94	The Wine Crate	
150	69	Vinos del Mundo	

ProductKey	ProdNr	ProdName	...
25	0178	Meerdael, Methode Traditionnelle Chardonnay, 2014	...
30	0199	Jacques Selosse, Brut Initial, 2012	
50	0212	Billecart-Salmon, Brut Réserve, 2014	

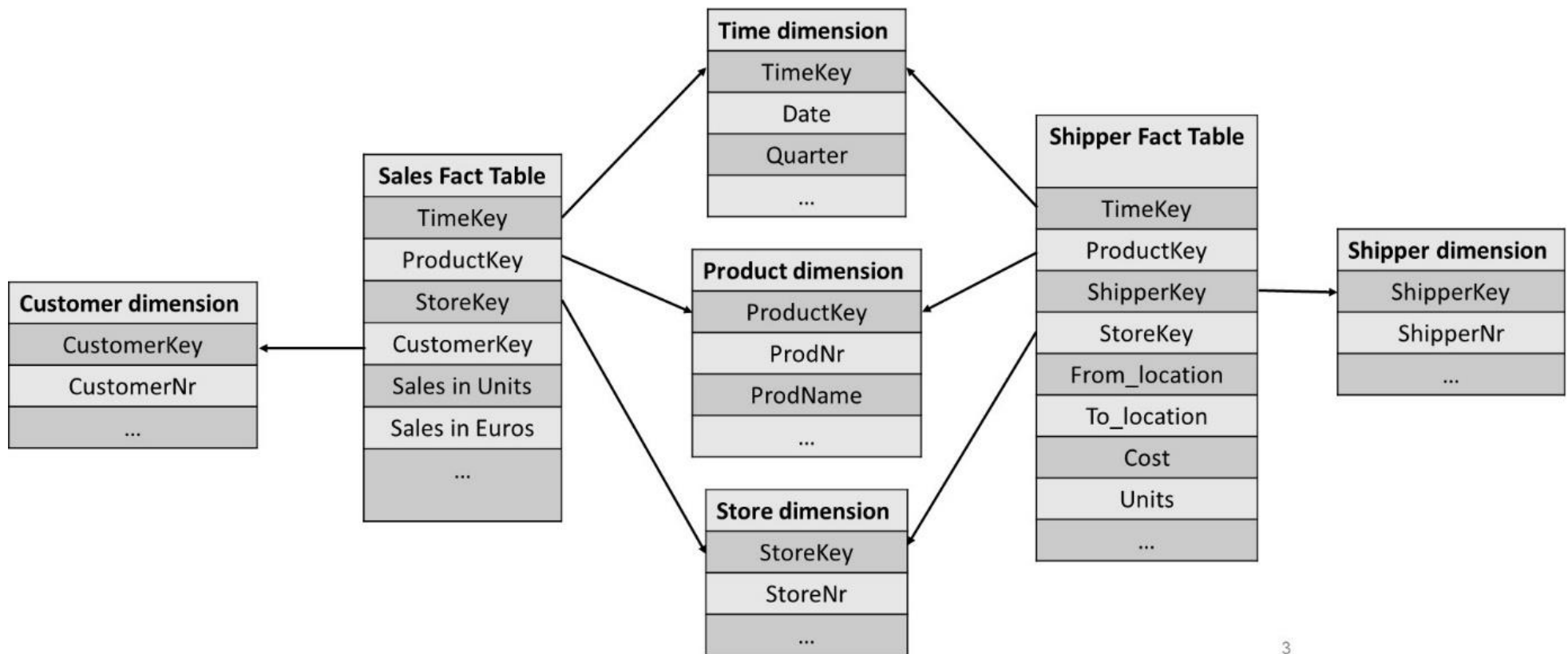
Snowflake Schema



Snowflake schema

- In case dimension tables grow too large
- In case most queries don't make use of the outer level dimension tables

Fact Constellation



3

Specific Schema Issues

- Surrogate keys
- Granularity of the fact table
- Factless Fact Tables
- Optimizing the dimension tables
- Defining junk dimensions
- Defining outrigger tables
- Slowly changing dimensions
- Rapidly changing dimensions

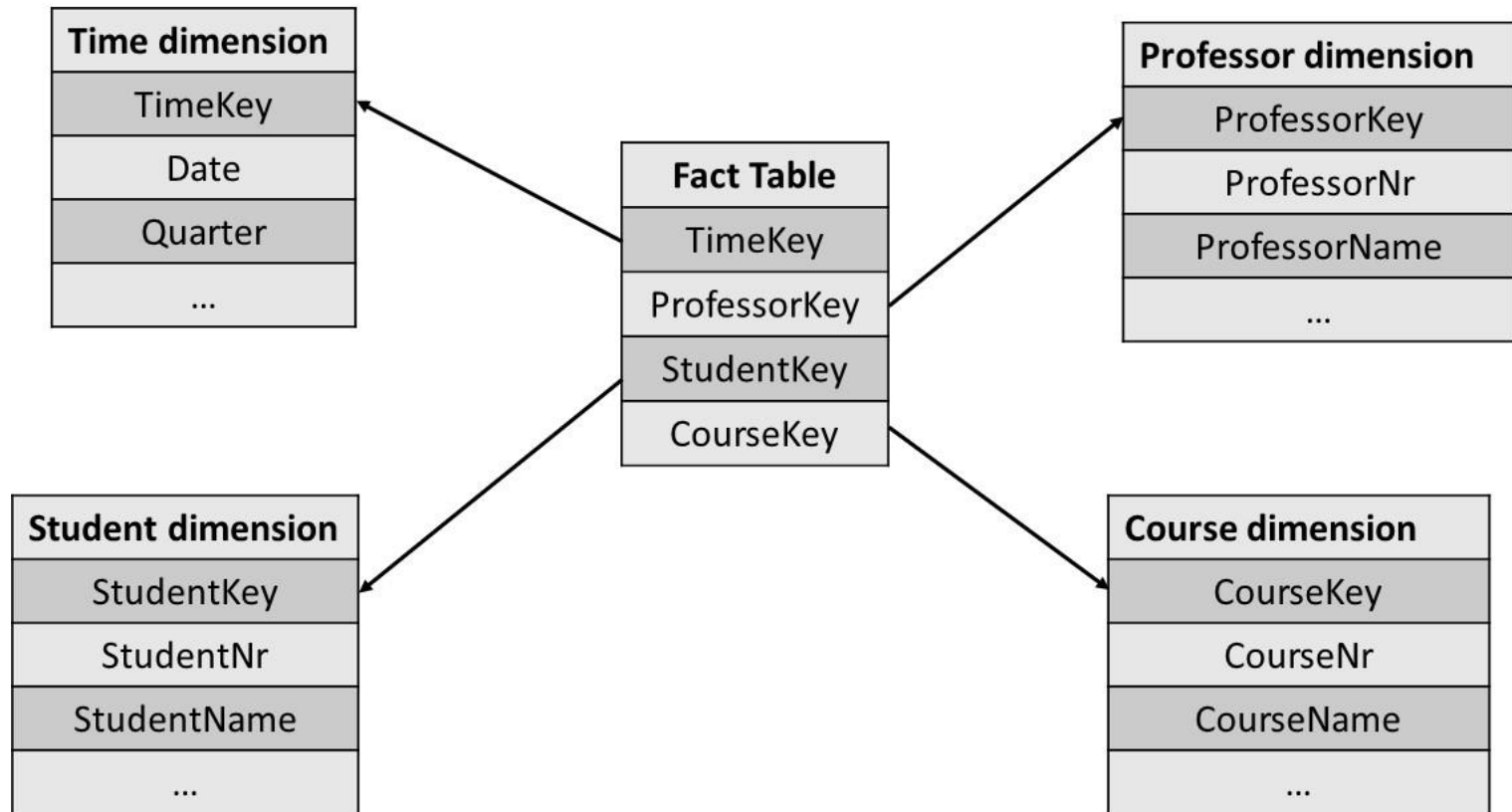
Surrogate keys

- StoreKey, ProductKey, ShipperKey, ...
- Meaningless integers
- Cannot use business keys since they usually have a business meaning
- Surrogate keys essentially buffer the data warehouse from the operational environment
- Business keys are usually bigger in size
- Business keys are also often re-used over longer periods of time

Granularity of the fact table

- Level of detail of one row of the fact table
- Higher (lower) granularity implies more (fewer) rows
- Trade-off between level of detailed analysis and storage requirements
- Examples:
 - One tuple of the fact table corresponds to one line on a purchase order
 - One tuple of the fact table corresponds to one purchase order
 - One tuple of the fact table corresponds to all purchase orders made by a customer

Factless Fact Tables



Optimizing the dimension tables

- Dimension tables should be heavily indexed to improve query execution time
- On average between 5 and 10
- E.g. Time
 - TimeKey, Date, DayOfWeek, DayOfMonth, DayOfYear, Month, MonthName, Year, LastDayInWeekFlag, LastDayInMonthFlag, FiscalWeek, HolidayFlag, ...

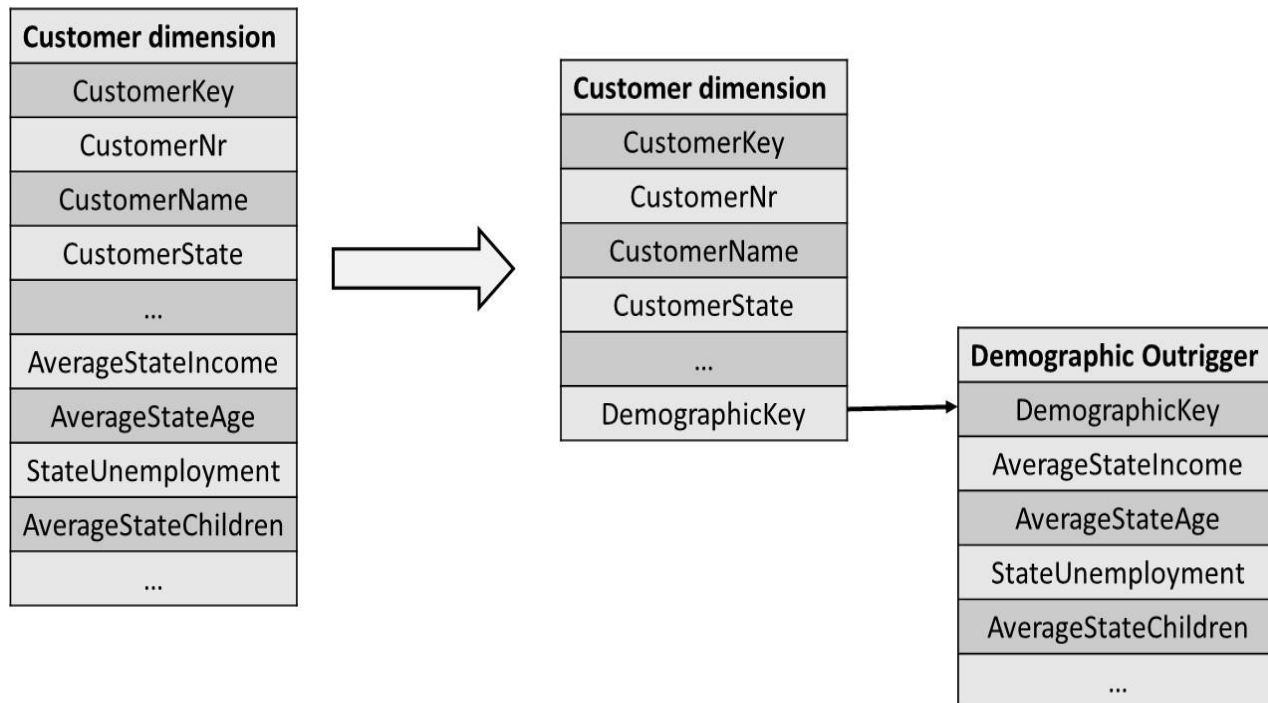
Junk Dimensions

- Deal with low cardinality attribute types such as flags or indicators
- Example: On-line Purchase (Y/N), Payment (cash or credit-card), Discount (Y/N)
- Junk dimension is a dimension that simply enumerates all feasible combinations of values of the low cardinality attribute types

JunkKey1	On-line purchase	Payment	Discount
1	Yes	Credit-card	Yes
2	Yes	Credit-card	No
3	No	Credit-card	Yes
4	No	Credit-card	No
5	No	Cash	Yes
6	No	Cash	No

Outtrigger tables

- store a set of attribute types of a dimension table which are highly correlated, low in cardinality and updated simultaneously



Slowly Changing Dimensions

- Dimensions that change slowly and irregularly over a period of time
- Example: customer segment ranging from AAA, AA, A, BBB, ... to C, determined on a yearly basis

Slowly Changing Dimensions

- **Approach 1**

OLD STATUS

CustomerKey	CustomerNr	CustomerName	Segment
123456	ABC123	Bart Baesens	AA

NEW STATUS

CustomerKey	CustomerNr	CustomerName	Segment
123456	ABC123	Bart Baesens	AAA

Slowly Changing Dimensions

- **Approach 2**

OLD STATUS

CustomerKey	CustomerNr	CustomerName	Segment	Start_Date	End_Date	Current_Flag
123456	ABC123	Bart Baesens	AA	27-02-2014	31-12-9999	Y

NEW STATUS

CustomerKey	CustomerNr	CustomerName	Segment	Start_Date	End_Date	Current_Flag
123456	ABC123	Bart Baesens	AA	27-02-2014	27-02-2015	N
123457	ABC123	Bart Baesens	AAA	28-02-2015	31-12-9999	Y

Slowly Changing Dimensions

- **Approach 3**

OLD STATUS

CustomerKey	CustomerNr	CustomerName	Segment
123456	ABC123	Bart Baesens	AA

NEW STATUS

CustomerKey	CustomerNr	CustomerName	Old Segment	New Segment
123456	ABC123	Bart Baesens	AA	AAA

Slowly Changing Dimensions

- **Approach 4**

OLD STATUS

CustomerKey	CustomerNr	CustomerName	Segment
123457	ABC123	Bart Baesens	AAA

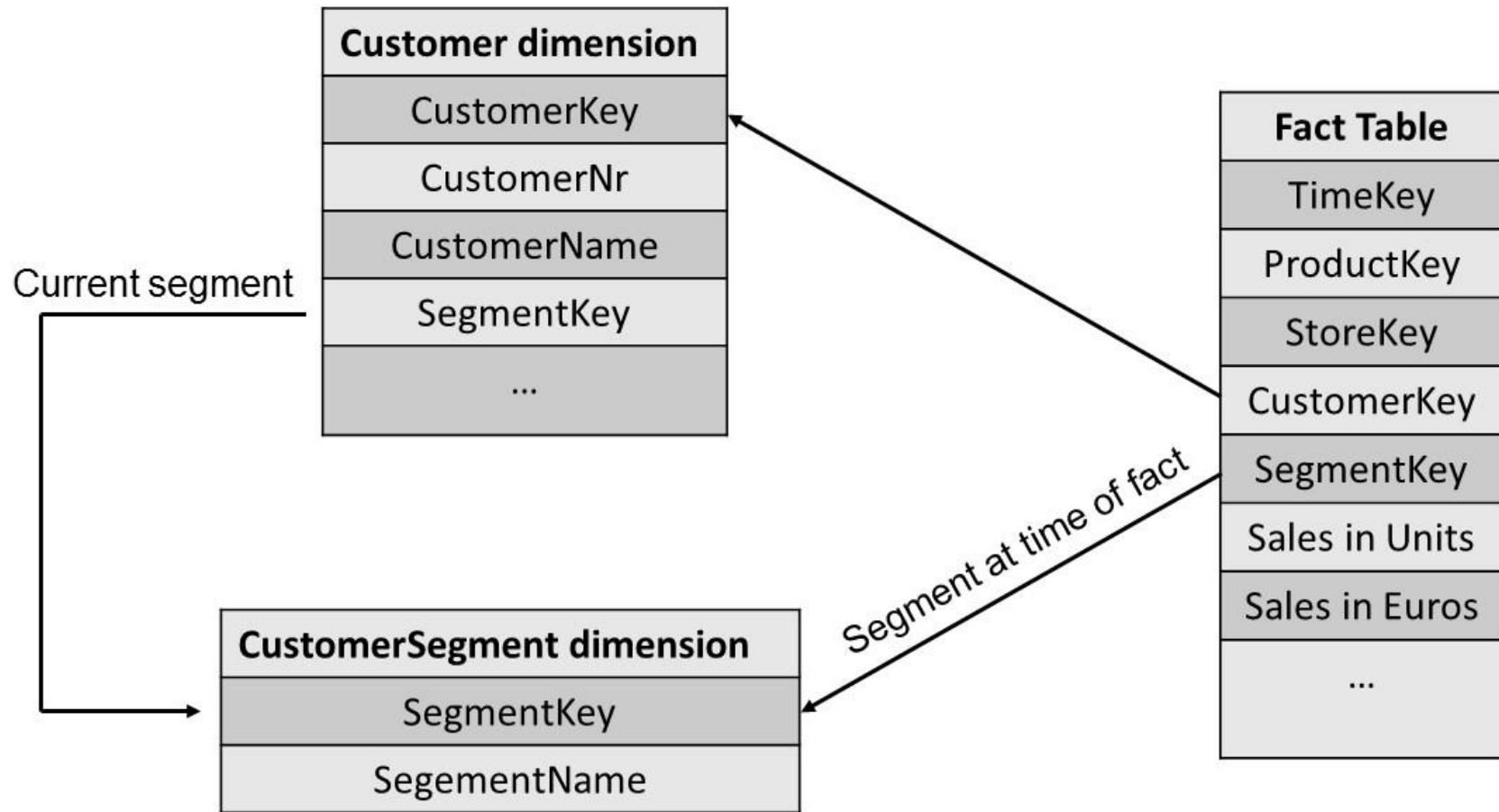
NEW STATUS

CustomerKey	CustomerNr	CustomerName	Segment	Start_Date	End_Date
123456	ABC123	Bart Baesens	AA	27-02-2014	27-02-2015
123457	ABC123	Bart Baesens	AAA	28-02-2015	31-12-9999

Rapidly changing dimensions

- Dimensions that change rapidly and regularly over a period of time
 - Example: customer segment ranging from AAA, AA, A, BBB, ... to C, determined on a daily basis
 - Approaches 2 and 4 discussed in the previous section will result into a lot of rows
 - Split customer information into stable (e.g., gender, marital status, ...) and rapidly changing information which is put in mini-dimension table
-

Rapidly changing dimensions



Rapidly changing dimensions

Fact table

TimeKey	ProductKey	StoreKey	CustomerKey	SegmentKey	Sales in Units	Sales in Euros
200	50	100	1200	1	6	167.94
210	25	130	1400	4	3	54
180	30	150	1000	3	12	384
...						

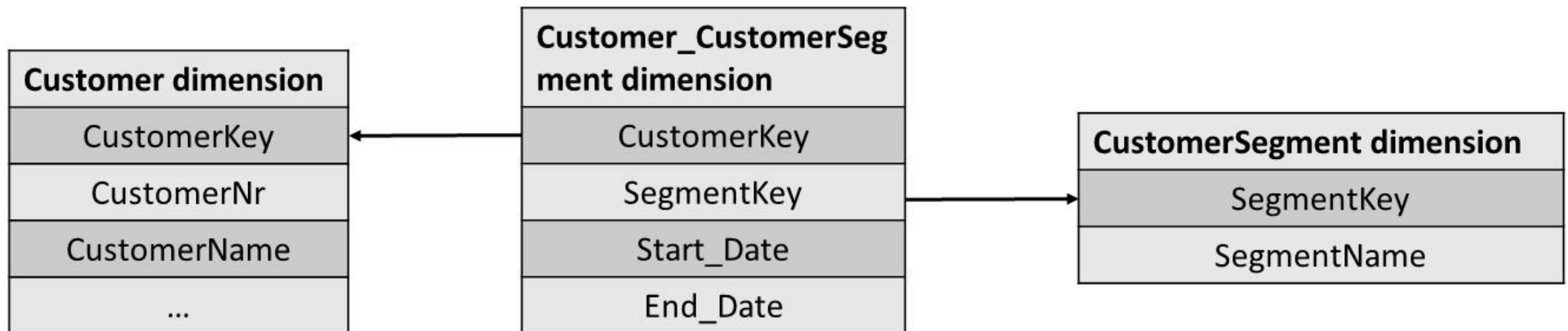
Customer Dimension

CustomerKey	CustomerNr	CustName	SegmentKey
1000	20	Bart Baesens	2
1200	10	Wilfried Lemahieu	1
1400	5	Seppe vanden Broucke	1

CustomerSegment Dimension

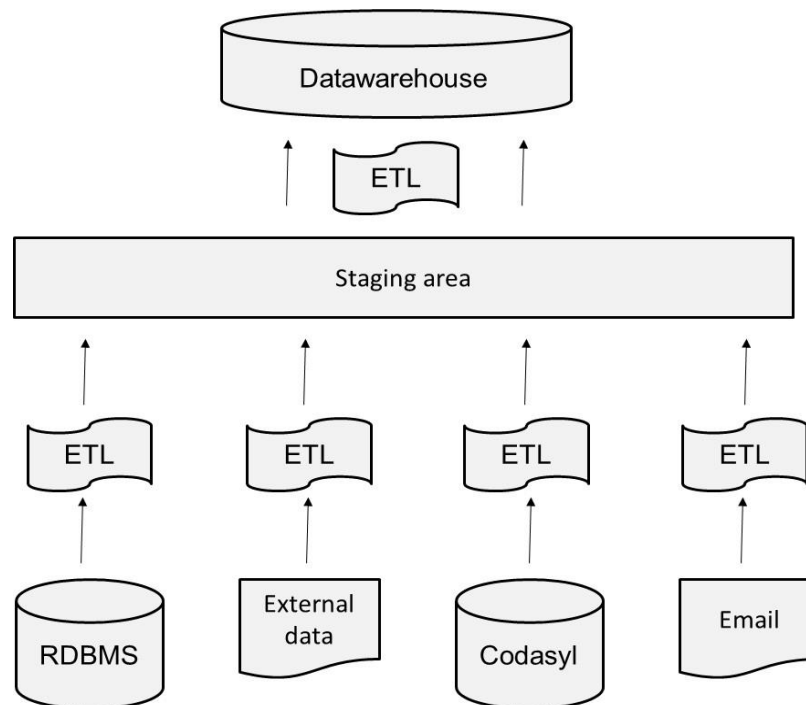
SegmentKey	SegmentName
1	A
2	B
3	C
4	D

Rapidly changing dimensions



The Extraction Transformation and Loading (ETL) Process

- can consume up to 80% of all efforts needed to set up a data warehouse
- dump the data in a so-called staging area where all the ETL activities can be executed



The Extraction Transformation and Loading (ETL) Process

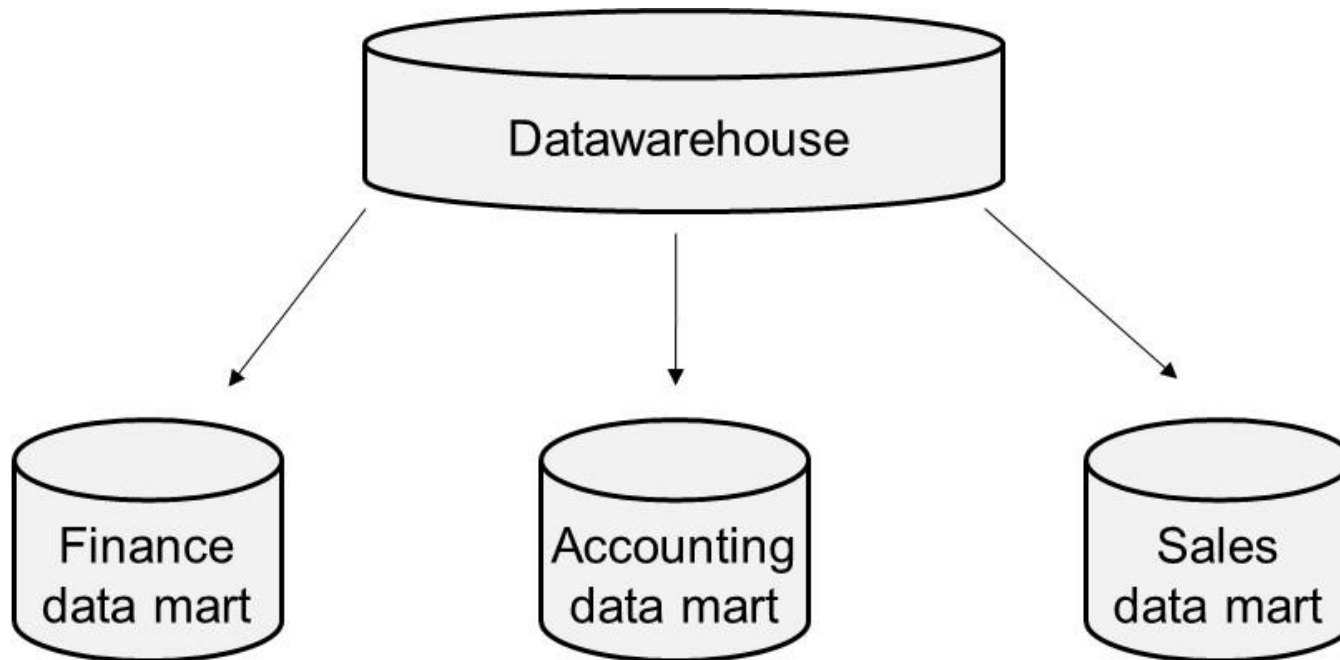
- Extraction
 - Full or incremental (Changed Data Capture)
- Transformation
 - formatting
 - cleansing
 - aggregation and merging
 - enrichment
- Loading
 - Fill the fact and dimension tables
- Documentation and metadata

Data Marts

- A **data mart** is a scaled down version of a data warehouse aimed at meeting the information needs of a homogeneous small group of end-users such as a department or business unit
- Provide focused content + improve query performance

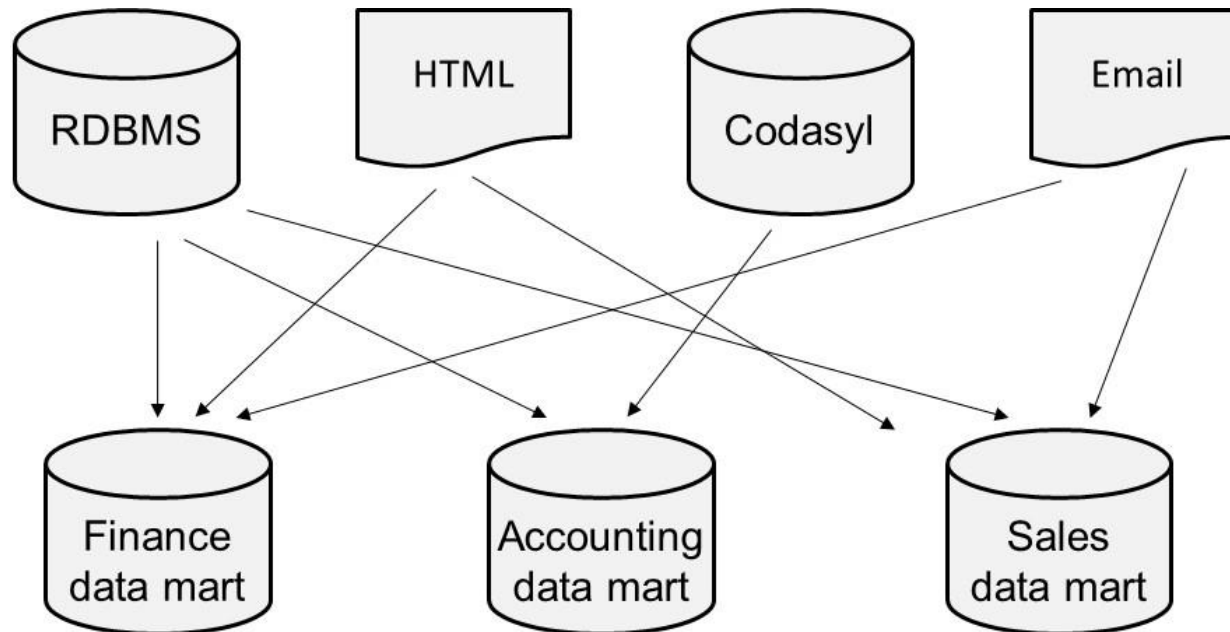
Data Marts

- **Dependent data marts** pull their data directly from a central data warehouse



Data Marts

- **Independent data marts** are standalone systems drawing data directly from the operational systems, external sources or a combination of both

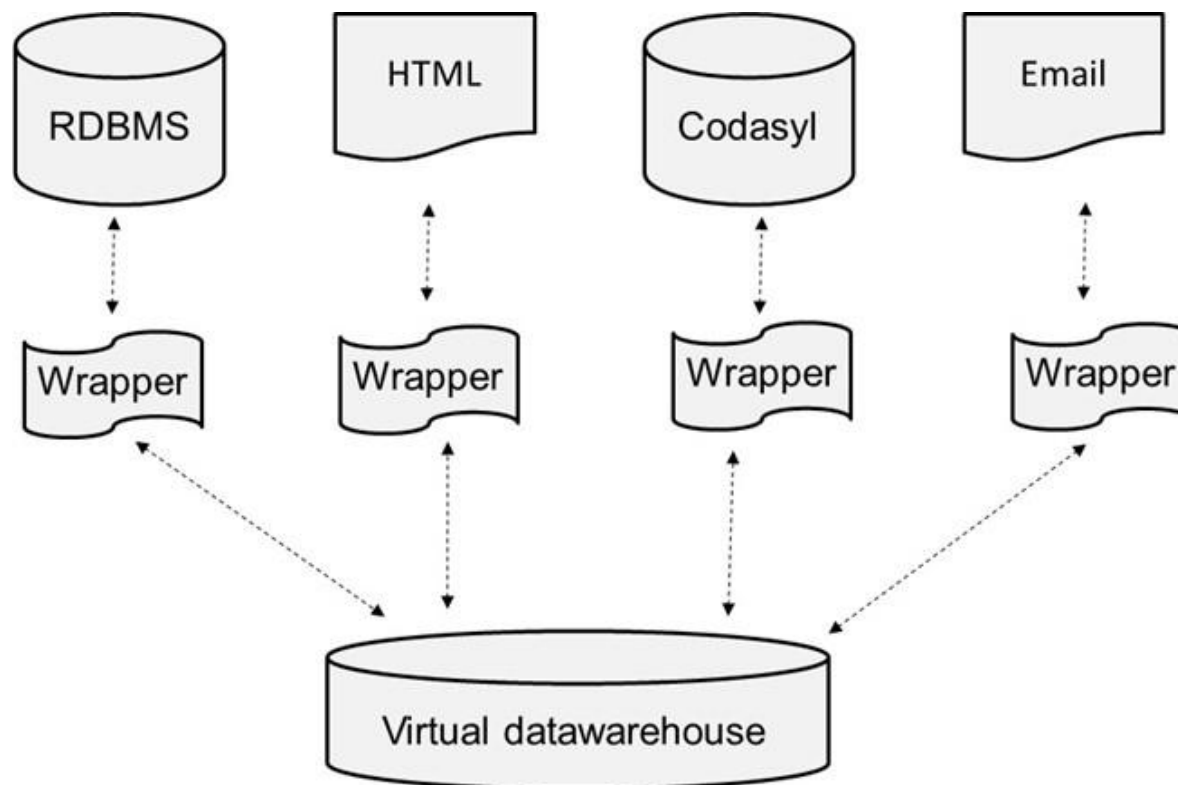


Virtual Data Warehouses and Virtual Data Marts

- Idea of virtualization is to use middleware to create a logical or **virtual data warehouse** or **virtual data mart** which has no physical data but provides a single point of access to a set of underlying physical data stores.
- Data is only accessed ('pulled') at query time.

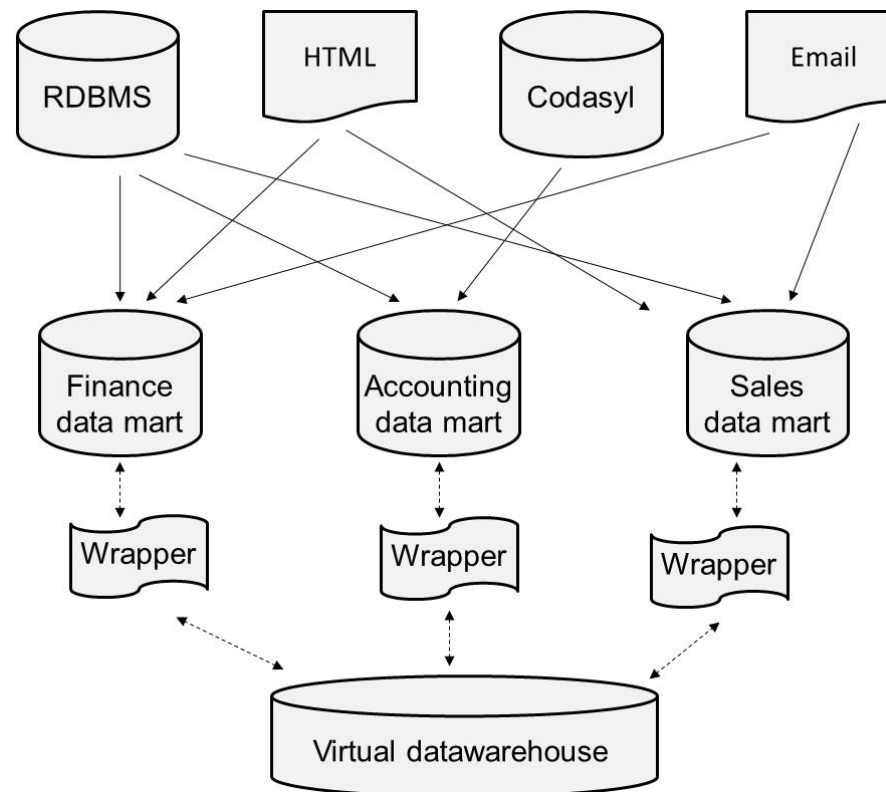
Virtual Data Warehouses and Virtual Data Marts

- Virtual Data warehouse can be built as a set of SQL views directly on the underlying operational data sources



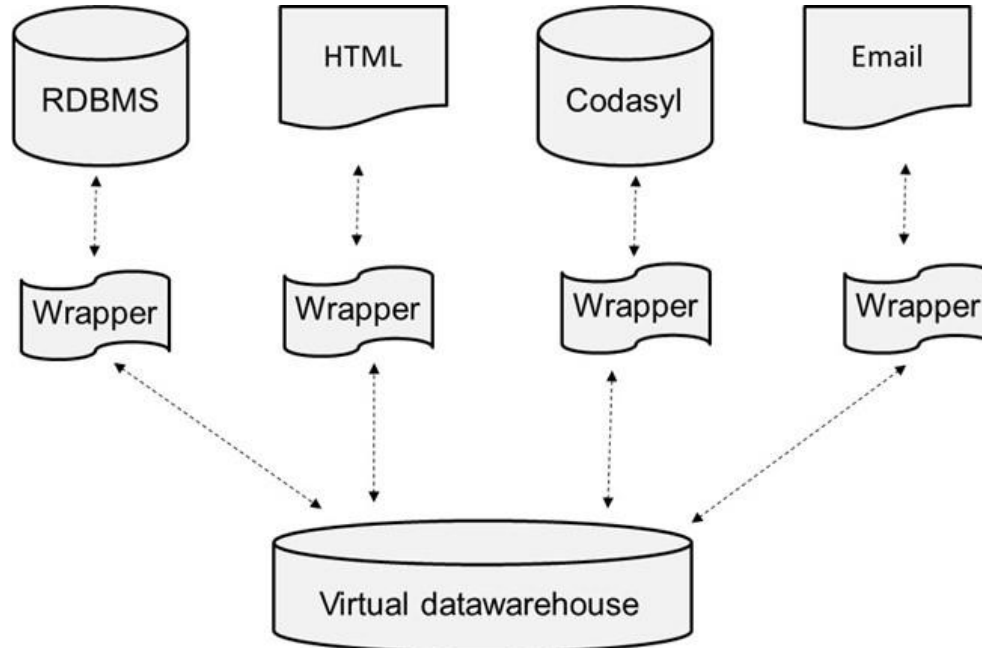
Virtual Data Warehouses and Virtual Data Marts

- Virtual data warehouse can be built as an extra layer on top of a collection of physical independent data marts



Virtual Data Warehouses and Virtual Data Marts

- Metadata model contains the schema mappings between the schemas of the underlying data stores and the schema of the virtual data warehouse
- Query reformulation



Virtual Data Warehouses and Virtual Data Marts

- A virtual data mart is usually defined as a single SQL view
- Virtual independent versus Virtual dependent data mart
- Disadvantages
 - Extra processing capacity from the underlying (operational) data sources
 - Not possible to keep track of historical data

Operational Data Store

- Operational Data Store (ODS) can be considered as a staging area that provides query facilities
- Good for analysis tools that need data that is closer to real time
- More complex analyses are still conducted on the actual data warehouse

Data Warehouses versus Data Lakes

	Data Warehouse	Data lake
Data	Structured	Often unstructured
Processing	Schema-on-write	Schema-on-read
Storage	Expensive	Low cost
Transformation	Before entering the DW	Before analysis
Agility	Low	High
Security	Mature	Maturing
Users	Decision makers	Data Scientists

Business Intelligence

- **Business intelligence (BI)** is referred to as the set of activities, techniques and tools aimed at understanding patterns in past data and predicting the future.
- Garbage In, Garbage Out (GIGO)
- BI techniques
 - Query and Reporting
 - Pivot tables
 - OLAP

Query and Reporting

- Business user can graphically and interactively design a query and corresponding report
- Self Service BI
- Query by Example (QBE)
 - a query is composed in a user-friendly and visual way
- Report can be refreshed at any time
- Innovative visualization techniques

Pivot tables

- A **pivot or cross-table** is a popular data summarization tool. It essentially cross-tabulates a set of dimensions in such a way that multidimensional data can be represented in a two-dimensional tabular format.

Sales		Quarter				<u>Total</u>
		Q1	Q2	Q3	Q4	
Region	Europe	100	200	50	100	450
	Africa	50	100	200	50	400
	Asia	20	50	10	150	230
	America	50	10	100	100	260
	<u>Total</u>	220	360	360	400	1340

On-Line Analytical Processing (OLAP)

- OLAP allows you interactively analyze the data, summarize it and visualize it in various ways
- Provide the business-user with a powerful tool for ad-hoc querying
- Types
 - MOLAP
 - ROLAP
 - HOLAP

MOLAP

- Multidimensional OLAP (MOLAP) stores the multidimensional data using a Multidimensional DBMS (MDBMS) whereby the data is stored in a multi-dimensional array-based data structure optimized for efficient storage and quick access.

<u>Array (key, value)</u>	Q1	Q2	Q3	Q4	<u>Total</u>
Product A	(1,1) 10	(1,2) 20	(1,3) 40	(1,4) 10	(1,5) 80
Product B	(2,1) 20	(2,2) 40	(2,3) 10	(2,4) 30	(2,5) 100
Product C	(3,1) 50	(3,2) 20	(3,3) 40	(3,4) 30	(3,5) 140
Product D	(4,1) 10	(4,2) 30	(4,3) 20	(4,4) 20	(4,5) 80
<u>Total</u>	(5,1) 90	(5,2) 110	(5,3) 110	(5,4) 90	(5,5) 400

MOLAP

- Fast in terms of data retrieval
- More storage space needed
- Scales poorly when the number of dimensions increases
- No universal SQL-like standard is provided
- Not optimized for transaction processing

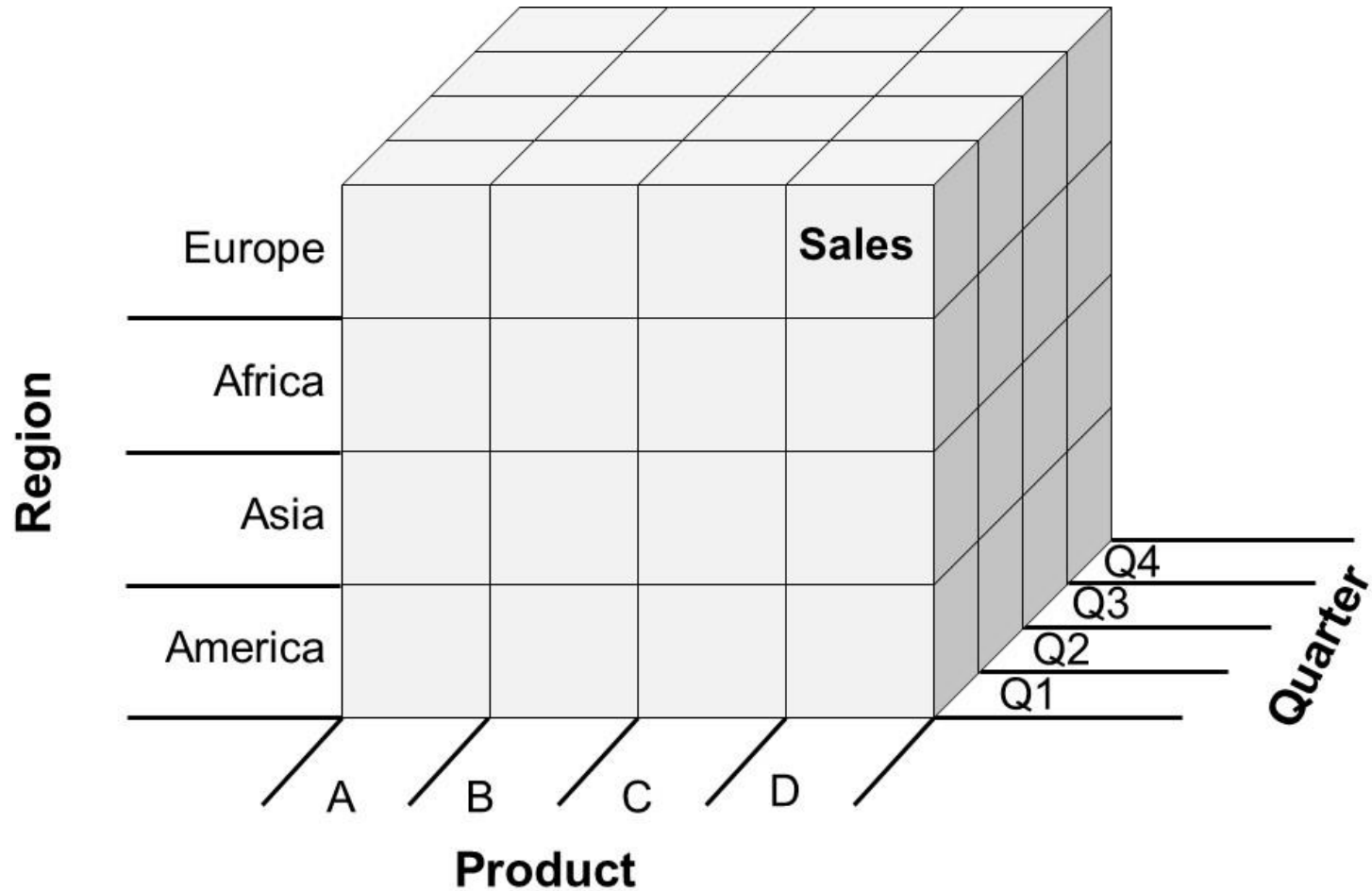
ROLAP

- **Relational OLAP (ROLAP)** stores the data in a relational data warehouse, which can be implemented using a star, snowflake or fact constellation schema
- ROLAP scales better to more dimensions than MOLAP
- ROLAP query performance may however be inferior to MOLAP

HOLAP

- **Hybrid OLAP (HOLAP)** tries to combine the best of both MOLAP and ROLAP
- RDBMS used to store the detailed data in a relational data warehouse whereas the pre-computed aggregated data is kept as a multidimensional array managed by an MDBMS
- OLAP analysis first starts from the multidimensional database
- Combine the performance of MOLAP with the scalability of ROLAP

OLAP operators



OLAP operators

- Roll-up refers to aggregating the current set of fact values within or across one or more dimensions
- Hierarchical versus dimensional roll-up

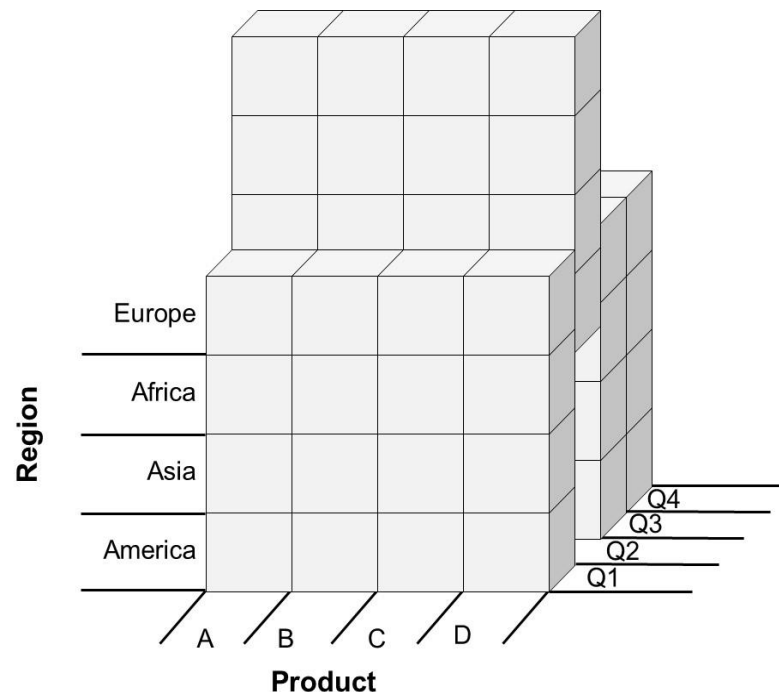
Region				
	Europe			
	Africa			
	Asia			
	America			
		A	B	C
		Product		

A	B	C	D
Product			

Reverse:
drill-down!

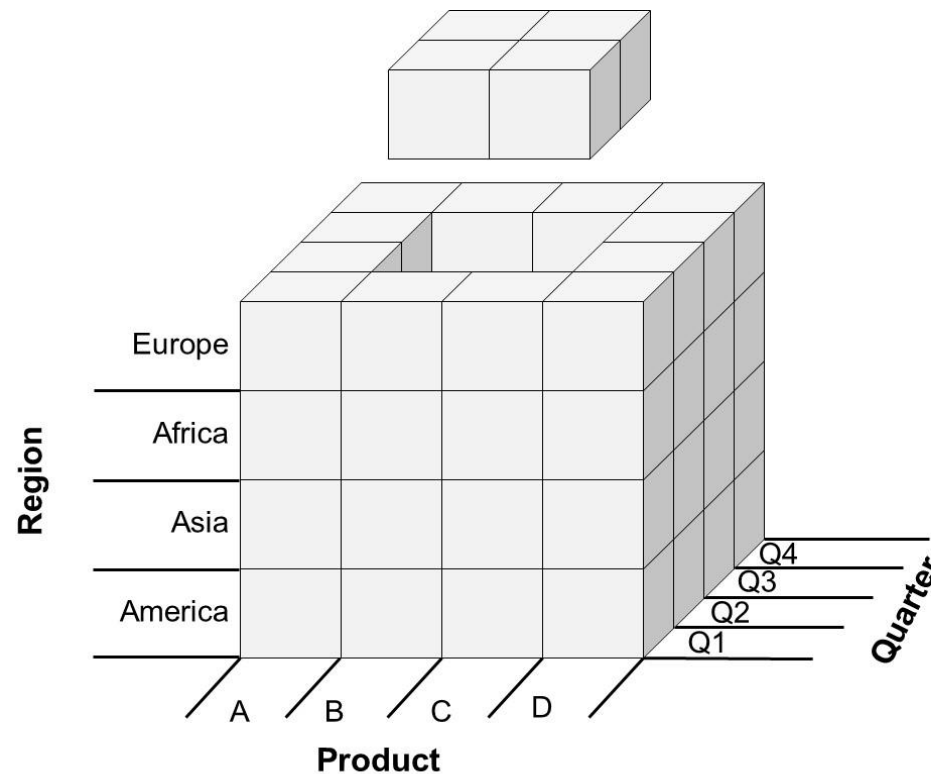
OLAP operators

- **Drill-across** is another OLAP operation whereby information from 2 or more connected fact tables is accessed.
- **Slicing** represents the operation whereby one of the dimensions is set at a particular value.



OLAP operators

- **Dicing** corresponds to a range selection on one or more dimensions



OLAP queries in SQL

- SQL-99 introduced 3 extensions to the GROUP BY statement: the CUBE, ROLLUP and GROUPING SETS operator
- The **CUBE** operator computes a union of GROUP BY's on every subset of the specified attribute types

OLAP queries in SQL

PRODUCT	QUARTER	REGION	SALES
A	Q1	Europe	10
A	Q1	America	20
A	Q2	Europe	20
A	Q2	America	50
A	Q3	America	20
A	Q4	Europe	10
A	Q4	America	30
B	Q1	Europe	40
B	Q1	America	60
B	Q2	Europe	20
B	Q2	America	10
B	Q3	America	20
B	Q4	Europe	10
B	Q4	America	40

OLAP queries in SQL

```
SELECT QUARTER, REGION, SUM(SALES)  
FROM SALESTABLE  
GROUP BY CUBE (QUARTER, REGION)
```

- this query computes the union of $2^2=4$ groupings of the SALESTABLE being: {(quarter,region), (quarter), (region), ()}, where () denotes an empty group list
- resulting multiset will have $4*2+4*1+1*2+1$ or 15 tuples

OLAP queries in SQL

QUARTER	REGION	SALES
Q1	Europe	50
Q1	America	80
Q2	Europe	40
Q2	America	60
Q3	Europe	NULL
Q3	America	40
Q4	Europe	20
Q4	America	80
Q1	NULL	130
Q2	NULL	100
Q3	NULL	40
Q4	NULL	90
NULL	Europe	110
NULL	America	250
NULL	NULL	360

CASE WHEN grouping(QUARTER) = 1
THEN 'All' ELSE QUARTER END AS QUARTER and
CASE WHEN grouping(REGION) = 1 THEN
'All' ELSE REGION END AS REGION

OLAP queries in SQL

- **ROLLUP** operator computes the union on every prefix of the list of specified attribute types, from the most detailed up to the grand total
- Key difference between the ROLLUP and CUBE operator is that the former generates a result set showing the aggregates for a hierarchy of values of the specified attribute types, whereas the latter generates a result set showing the aggregates for all combinations of values of the selected attribute types

OLAP queries in SQL

```
SELECT QUARTER, REGION, SUM(SALES)  
FROM SALESTABLE  
GROUP BY ROLLUP (QUARTER, REGION)
```

- This query generates the union of three groupings $\{(\text{quarter}, \text{region}), (\text{quarter}), ()\}$ where $()$ again represents the full aggregation
- Resulting multiset will thus have $4*2+4+1$ or 13 rows

OLAP queries in SQL

QUARTER	REGION	SALES
Q1	Europe	50
Q1	America	80
Q2	Europe	40
Q2	America	60
Q3	Europe	NULL
Q3	America	40
Q4	Europe	20
Q4	America	80
Q1	NULL	130
Q2	NULL	100
Q3	NULL	40
Q4	NULL	90
NULL	NULL	360

OLAP queries in SQL

- GROUP BY ROLLUP construct can also be applied to attribute types that represent different aggregation levels along the same dimension

```
SELECT REGION, COUNTRY, CITY, SUM(SALES)  
FROM SALESTABLE  
GROUP BY ROLLUP (REGION, COUNTRY, CITY)
```

OLAP queries in SQL

- **GROUPING SETS** operator generates a result set equivalent to that generated by a UNION ALL of multiple simple GROUP BY clauses

```
SELECT QUARTER, REGION, SUM(SALES)
FROM SALESTABLE
GROUP BY GROUPING SETS ((QUARTER), (REGION))
```

```
SELECT QUARTER, NULL, SUM(SALES)
FROM SALESTABLE
GROUP BY QUARTER
UNION ALL
SELECT NULL, REGION, SUM(SALES)
FROM SALESTABLE
GROUP BY REGION
```

OLAP Queries in SQL

QUARTER	REGION	SALES
Q1	NULL	130
Q2	NULL	100
Q3	NULL	40
Q4	NULL	90
NULL	Europe	110
NULL	America	250

OLAP Queries in SQL

```
SELECT QUARTER, REGION, SUM(SALES)  
FROM SALESTABLE  
GROUP BY CUBE (QUARTER, REGION)
```

```
SELECT QUARTER, REGION, SUM(SALES)  
FROM SALESTABLE  
GROUP BY GROUPING SETS ((QUARTER,  
REGION), (QUARTER), (REGION), ( ))
```

OLAP Queries in SQL

```
SELECT QUARTER, REGION, SUM(SALES)  
FROM SALESTABLE  
GROUP BY ROLLUP (QUARTER, REGION)
```

```
SELECT QUARTER, REGION, SUM(SALES)  
FROM SALESTABLE  
GROUP BY GROUPING SETS ((QUARTER,  
REGION), (QUARTER), ())
```

OLAP Queries in SQL

- SQL2003 introduced support for 2 types of other activities: ranking and windowing

OLAP Queries in SQL

PRODUCT	SALES
A	50
B	20
C	10
D	45
E	40
F	30
G	60
H	20
I	15
J	25

OLAP Queries in SQL

```
SELECT PRODUCT, SALES,  
RANK() OVER (ORDER BY SALES ASC) as  
RANK_SALES,  
DENSE_RANK() OVER (ORDER BY SALES ASC) as  
DENSE_RANK_SALES, PERCENT_RANK() OVER (ORDER  
BY SALES ASC) as PERC_RANK_SALES,  
CUM_DIST() OVER (ORDER BY SALES ASC) as  
CUM_DIST_SALES,  
FROM SALES  
ORDER BY RANK_SALES ASC
```

OLAP Queries in SQL

Product	Sales	RANK_SALES	DENSE_RANK_SALES	PERC_RANK_SALES	CUM_DIST_SALES
C	10	1	1	0	0.1
I	15	2	2	1/9=0.11	0.2
B	20	3	3	2/9=0,22	0.4
H	20	3	3	2/9=0,22	0.4
J	25	5	4	4/9=0,44	0.5
F	30	6	5	5/9=0,55	0.6
E	40	7	6	6/9=0,66	0.7
D	45	8	7	7/9=0,77	0.8
A	50	9	8	8/9=0,88	0.9
G	60	10	9	9/9=1	1

OLAP Queries in SQL

- **Windowing** allows calculating cumulative totals or running averages based on a specified time window.

QUARTER	REGION	SALES
1	America	10
2	America	20
3	America	10
4	America	30
1	Europe	10
2	Europe	20
3	Europe	10
4	Europe	20

OLAP Queries in SQL

```
SELECT QUARTER, REGION, SALES,  
AVG(SALES) OVER (PARTITION BY REGION  
ORDER BY QUARTER ROWS BETWEEN 1  
PRECEDING AND 1 FOLLOWING) AS SALES_AVG  
FROM SALES  
ORDER BY REGION, QUARTER, SALES_AVG
```


OLAP Queries in SQL

QUARTER	REGION	SALES	SALES_AVG
1	America	10	15
2	America	20	13,33
3	America	10	20
4	America	30	20
1	Europe	10	15
2	Europe	20	13,33
3	Europe	10	16,67
4	Europe	20	15

Conclusions

- Operational versus Tactical/Strategic Decision Making
- Data Warehouse Definition
- Data Warehouse Schemas
- The Extraction Transformation and Loading (ETL) Process
- Data Marts
- Virtual Data Warehouses and Virtual Data Marts
- Operational Data Store
- Data Warehouses versus Data Lakes
- Business intelligence