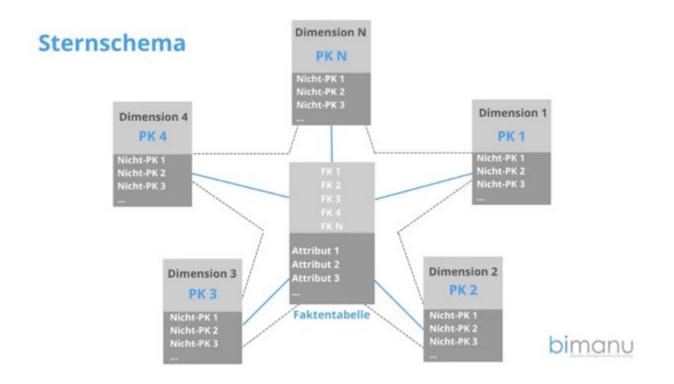


# Data Engineer

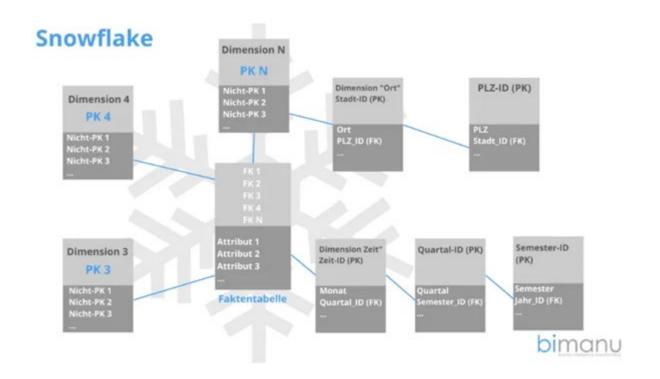
Snowflake - Schema:

Snowflake Schema Grundlagen, Datenmodellierung, Erstellung Snowflake Schema in RDBMS



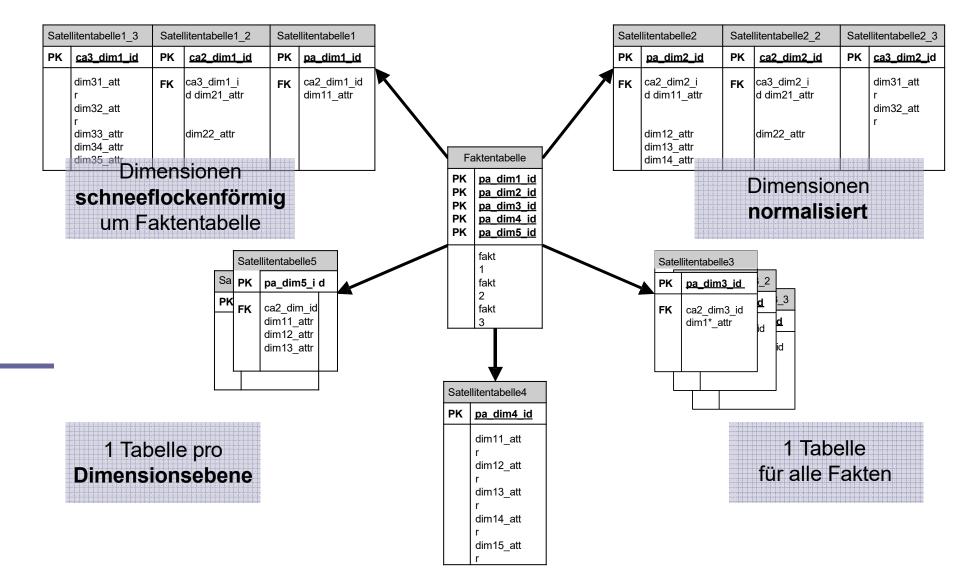






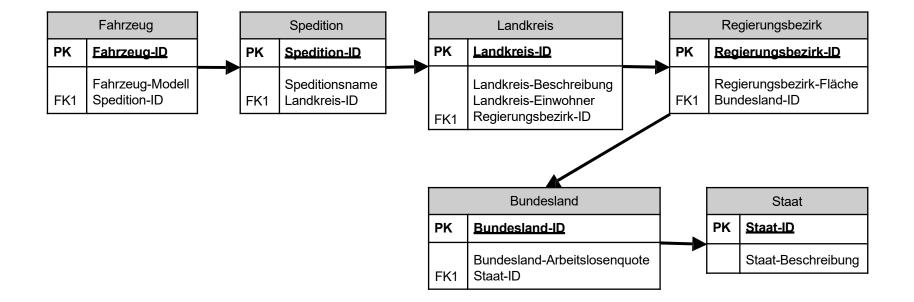
## Snowflake-Schema





# Snowflake-Dimension: Beispiel





### Snowflake-Schema



- Warum muss jetzt alles ID heißen?
  - Eindeutigkeit der Referenzkette in der Klassifkationshierarchie:
    - Eine Monats-ID aus der Faktentabelle gehört zu einem bestimmten Quartal in einem bestimmten Jahr
    - "2. Quartal" reicht dafür nicht, weil das Jahr nicht mehr herausgefunden werden kann
- → gegebenfalls sind zu den IDs in den Satellitentabellen noch Beschreibungen zu ergänzen (Speditionsname, Landkreis-Beschreibung, ...)

## Snowflake-Schema



- Probleme
  - weniger übersichtlich
  - weniger performante Anfragebearbeitung
- Eigenschaften
  - normalisierte Dimensionstabellen
  - Dimensionsklassifikationen explizit modelliert

### Star- vs. Snowflake-Schema



- Welche Variante ist besser?
  - → Welches Schema geeignet ist, hängt vom Anwendungsprofil ab!
  - → Dennoch wird sehr häufig ein Star-Schema oder ein leicht modifiziertes Star-Schema verwendet.

### **OLAP-** Architekturen



- **ROLAP:** Relational On Line Analytical Processing, relationale Datenspeicherung Tabellenform.
- MOLAP: Multidimensional On Line Analytical Processing, multidimensional Datenspeicherung, n-dimensionaler Würfel.
- **HOLAP:** Hybrid On Line Analytical Processing (HOLAP) Speicherung eines Teils des DWH's in Form von Würfeln.

#### **Data Warehouse**



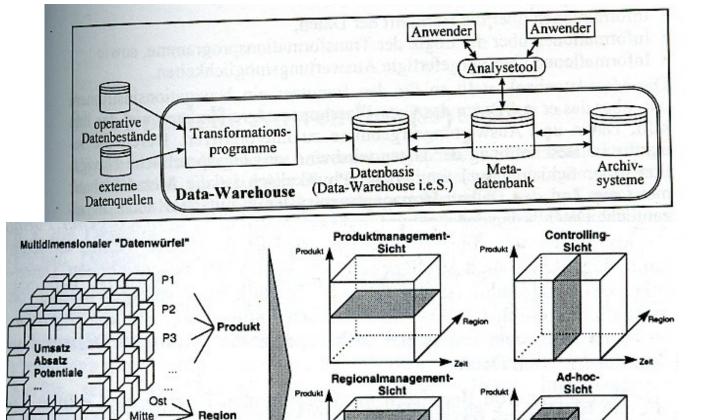
#### **Grundkonzept Datenbank/Datenmodellierung**

Jan. Feb. Mårz ...

Zeit

 Ziel: Bereitstellung aller relevanten Daten für betriebswirtschaftliche Analysen des Managements

- Ausprägungen
  - global = Data-Warehouse
  - spezifisch = Data-Mart
- Schritte
  - Modellierung
  - Data-Mining
  - Präsentation



Region

Region

# Multi-dimensional Schemas (2)



- Two common multi-dimensional schemas are
  - Star Schema:
    - Consists of a fact table with a single table for each dimension
  - Snowflake Schema:
    - It is a variation of Star Schema, in which the dimensional tables from a Snowflake Schema are organized into a hierarchy by normalizing them.

### Conceptual Modeling of Data Warehouses

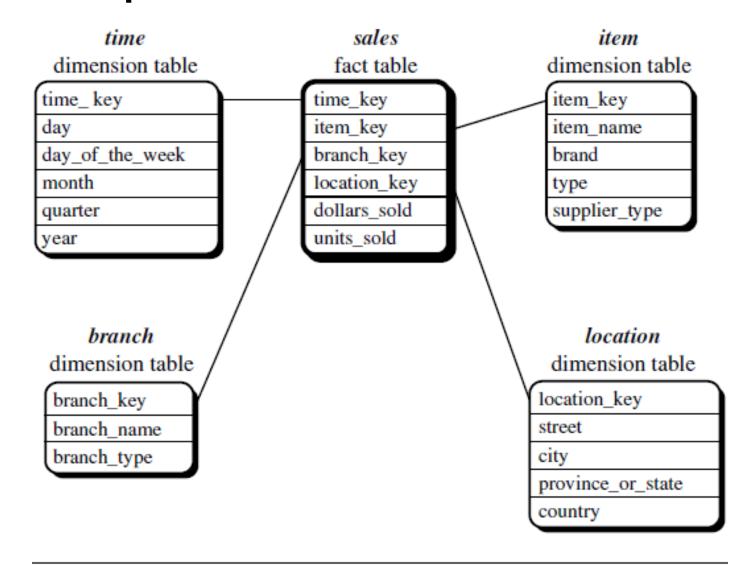


- Modeling data warehouses: multidimensional model
  - Star schema: A fact table in the middle connected to a set of dimension tables
  - Snowflake schema: A refinement of star schema where some dimensional hierarchy is normalized into a set of smaller dimension tables, forming a shape similar to snowflake

May 15, 2025 12

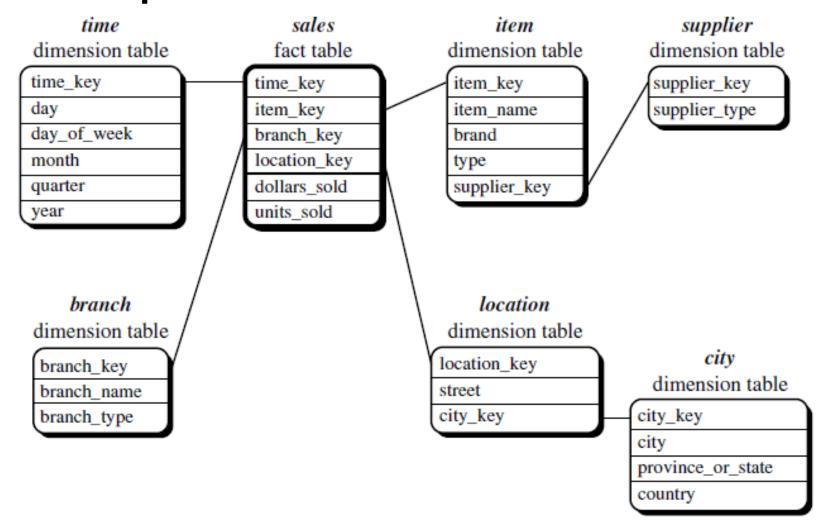
## Example of Star Schema





## Example of Snowflake Schema

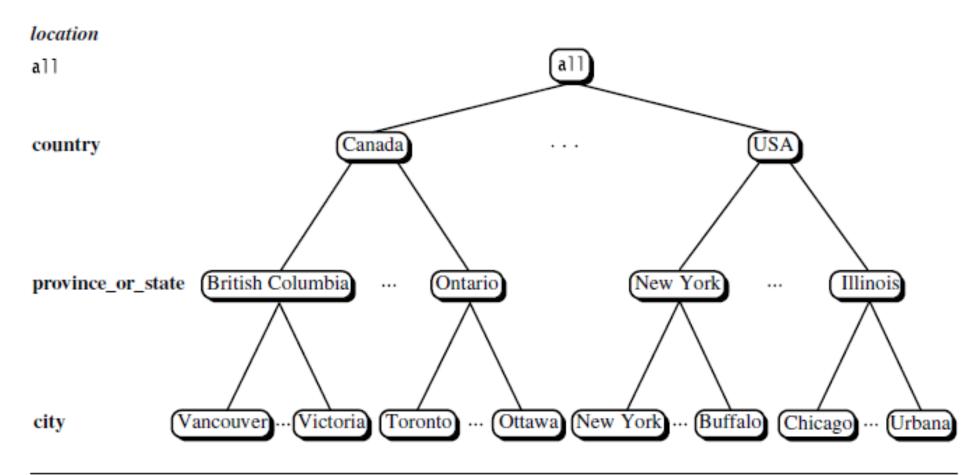




May 15, 2025 14

### A Concept Hierarchy: Dimension (location)



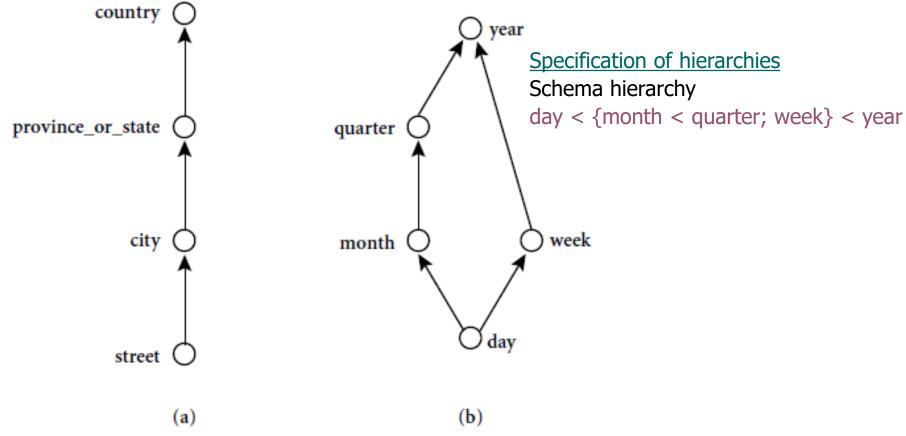


A concept hierarchy for the dimension location

May 15, 2025

### Concept Hierarchies





Hierarchical and lattice structures of attributes in warehouse dimensions:
(a) a hierarchy for *location*; (b) a lattice for *time* 

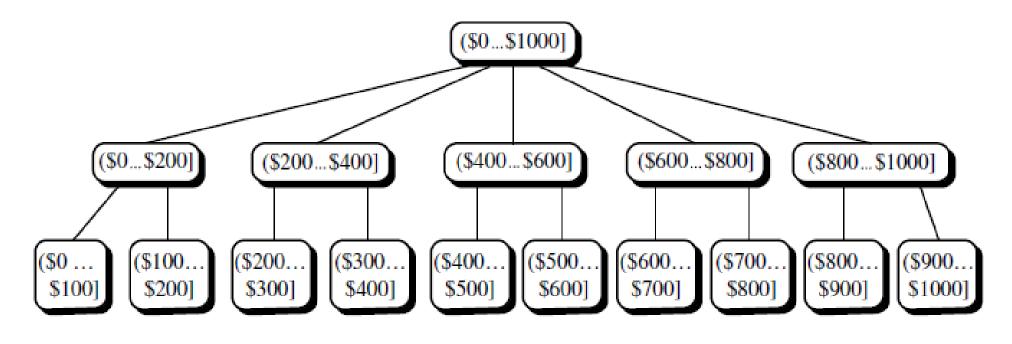
May 15, 2025

### Concept Hierarchies



### **Specification of hierarchies**

Set\_grouping hierarchy

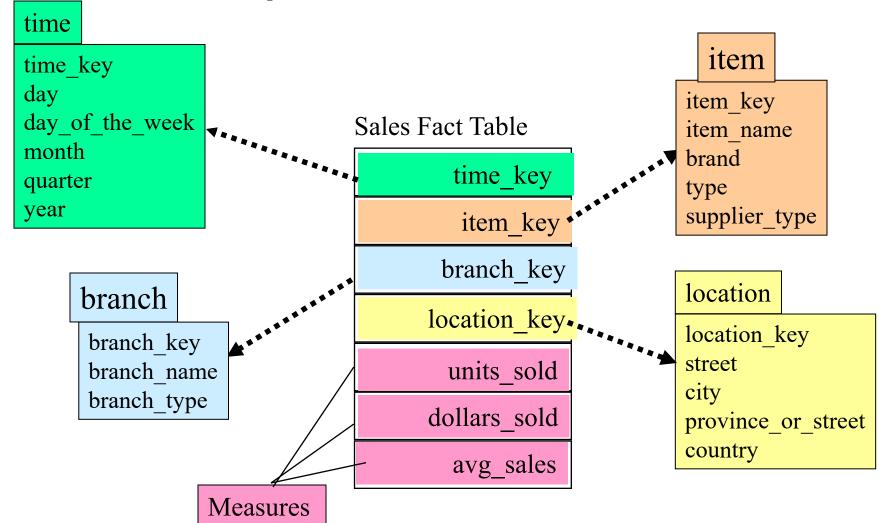


A concept hierarchy for the attribute price

May 15, 2025 17

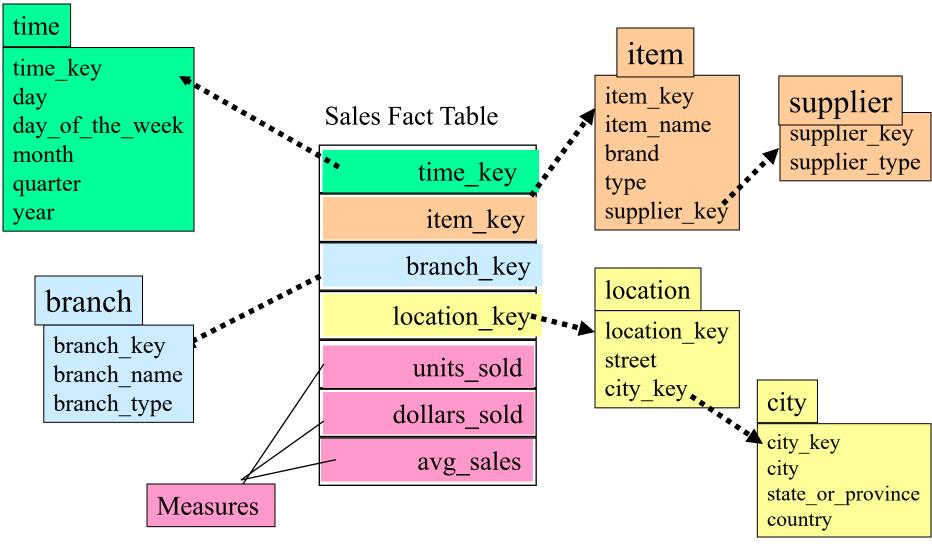


# Example of Star Schema



# Example of Snowflake Schema







## **Dimension Hierarchies**



store	storeld	cityld	tld	mgr
	s5	sfo	t1	joe
	s7	sfo	t2	fred
	s9	la	t1	nancy

sType	tld	size	location
	t1	small	downtown
	t2	large	suburbs

	city	cityld	pop	regld
`		sfo	1M	north
		la	5M	south

- → snowflake schema
- → constellations

region	regld	name
	north	cold region
	south	warm region

## Cube



### Fact table view:

# sale prodId storeId amt p1 c1 12 p2 c1 11 p1 c3 50 p2 c2 8

### Multi-dimensional cube:

		с1	c2	c3
<b>&gt;</b>	p1	12		50
	p2	11	8	

dimensions = 2

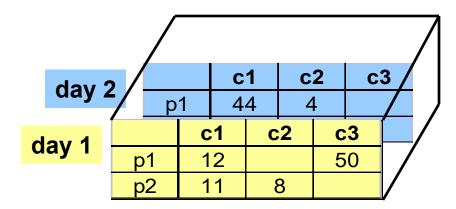
## 3-D Cube



### Fact table view:

sale	prodld	storeld	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	с3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

### Multi-dimensional cube:



dimensions = 3

# Aggregates



- Add up amounts for day 1
- In SQL: SELECT sum(amt) FROM SALE WHERE date = 1

sale	prodld	storeld	date	amt
	p1	с1	1	12
	p2	с1	1	11
	p1	сЗ	1	50
	p2 p1	c2	1	8
	p1	с1	2	44
	p1	c2	2	4



81

# Aggregates



- Add up amounts by day
- In SQL: SELECT date, sum(amt) FROM SALE GROUP BY date

sale	prodld	storeld	date	amt
	p1	c1	1	12
	p2	с1	1	11
	p1	с3	1	50
	p2	c2	1	8
	p1	с1	2	44
	p1	с2	2	4



ans	date	sum
	1	81
	2	48

# Another Example



- Add up amounts by day, product
- In SQL: SELECT date, sum(amt) FROM SALE GROUP BY date, prodld

sale	prodld	storeld	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	с3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4



sale	prodld	date	amt
	p1	1	62
	p2	1	19
	p1	2	48





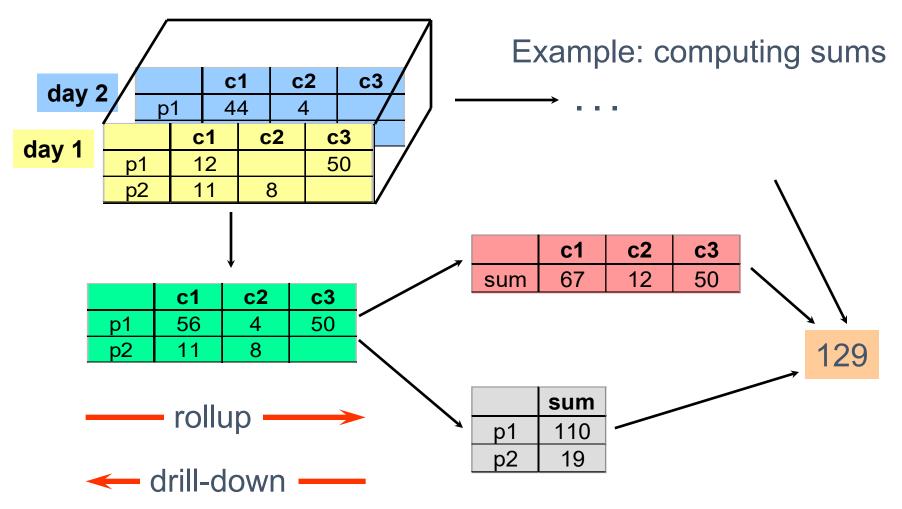
# Aggregates

- Operators: sum, count, max, min,
- "Having" clause
- Using dimension hierarchy
  - average by region (within store)
  - maximum by month (within date)

median, ave

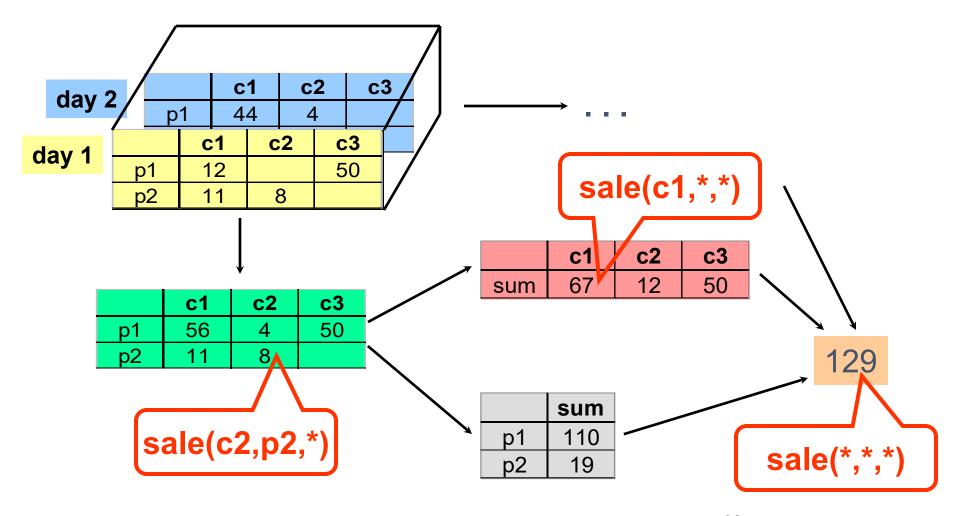
# **Cube Aggregation**





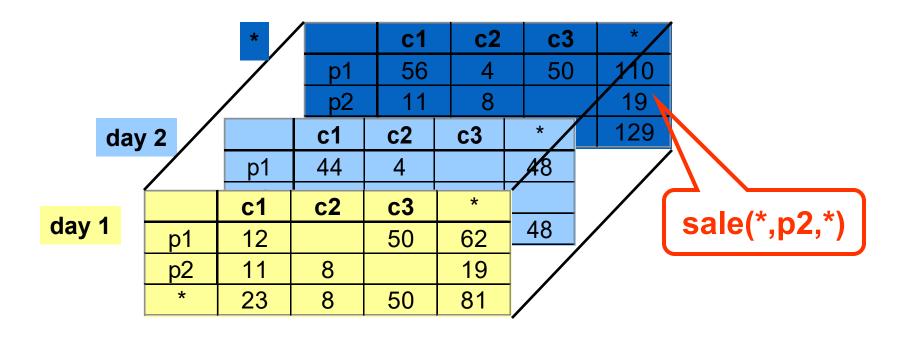
# **Cube Operators**





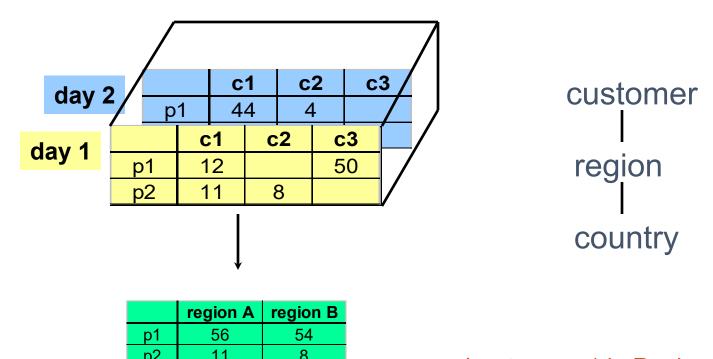
## **Extended Cube**





## Aggregation Using Hierarchies





(customer c1 in Region A; customers c2, c3 in Region B)

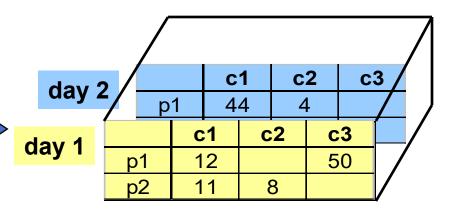
# Pivoting



### Fact table view:

sale	prodld	storeld	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

### Multi-dimensional cube:





	с1	c2	с3
p1	56	4	50
p2	11	8	

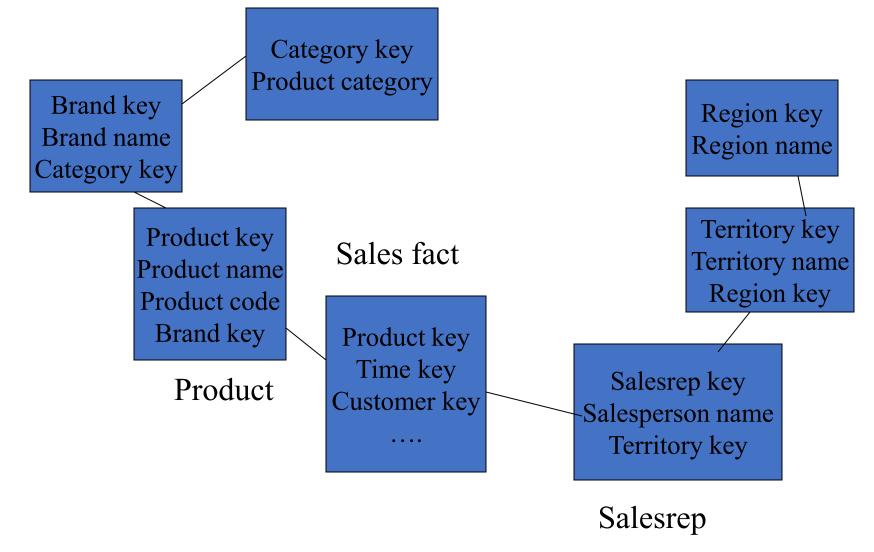




- Snowflake schema is a type of star schema but a more complex model.
- "Snowflaking" is a method of normalizing the dimension tables in a star schema.
- The normalization eliminates redundancy.
- The result is more complex queries and reduced query performance.

## Sales: Snowflake Schema

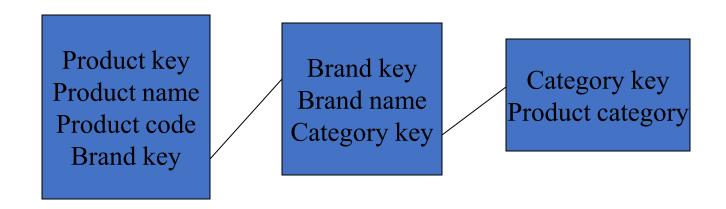




# Snowflaking



 The attributes with low cardinality in each original dimension table are removed to form separate tables. These new tables are linked back to the original dimension table through artificial keys.



### Snowflake Schema

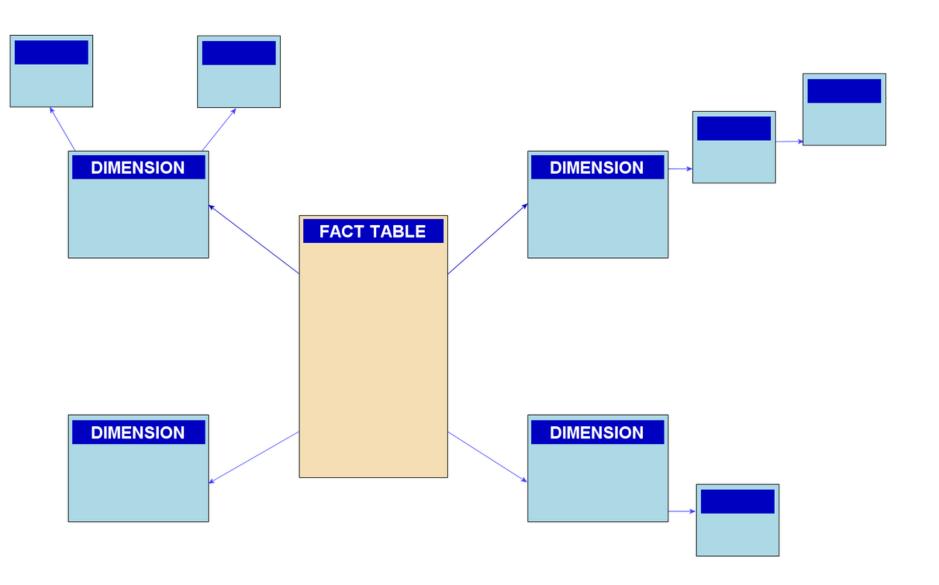


- Advantages:
  - Small saving in storage space
  - Normalized structures are easier to update and maintain
- Disadvantages:
  - Schema less intuitive and end-users are put off by the complexity
  - Ability to browse through the contents difficult
  - Degrade query performance because of additional joins



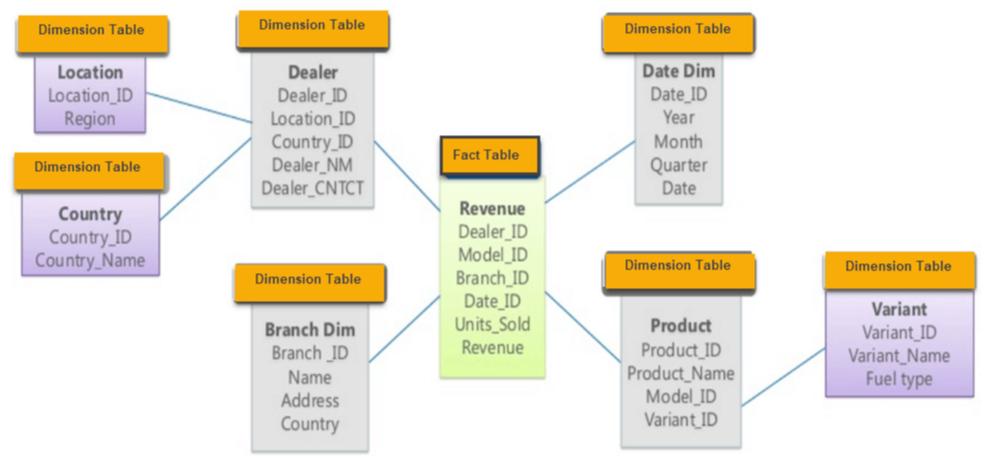


- Performance benchmarking can be used to determine what is the best design.
- Snowflake schema: easier to maintain dimension tables when dimension tables are very large (reduce overall space). It is not generally recommended in a data warehouse environment.
- Star schema: more effective for data cube browsing (less joins): can affect performance.

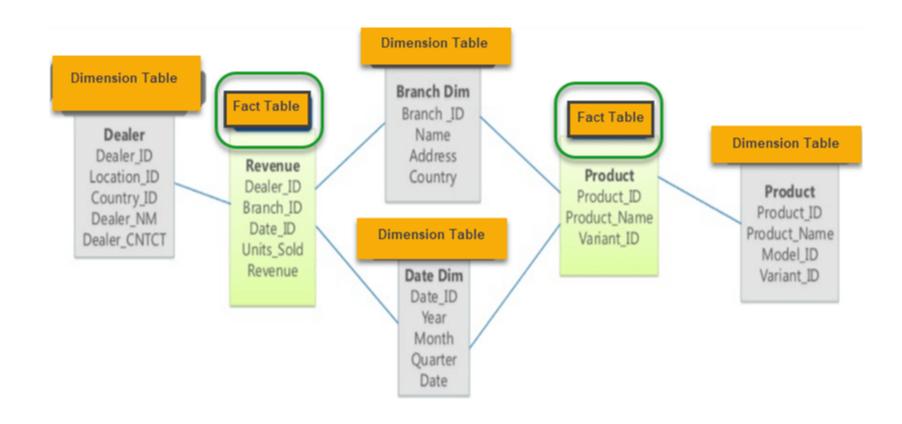




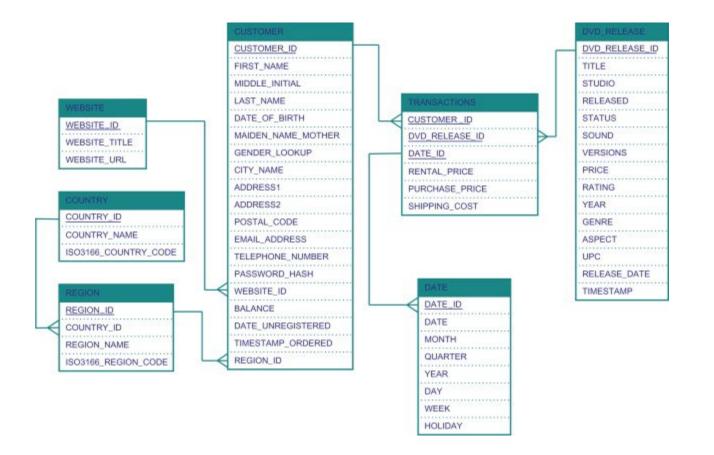




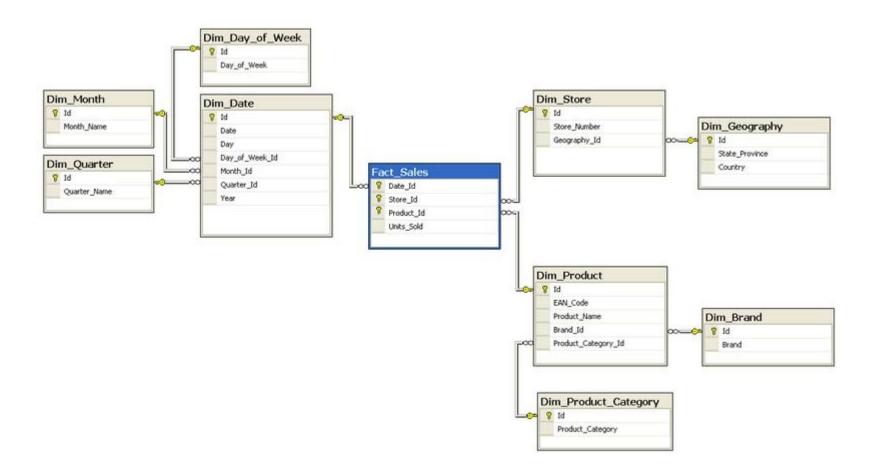


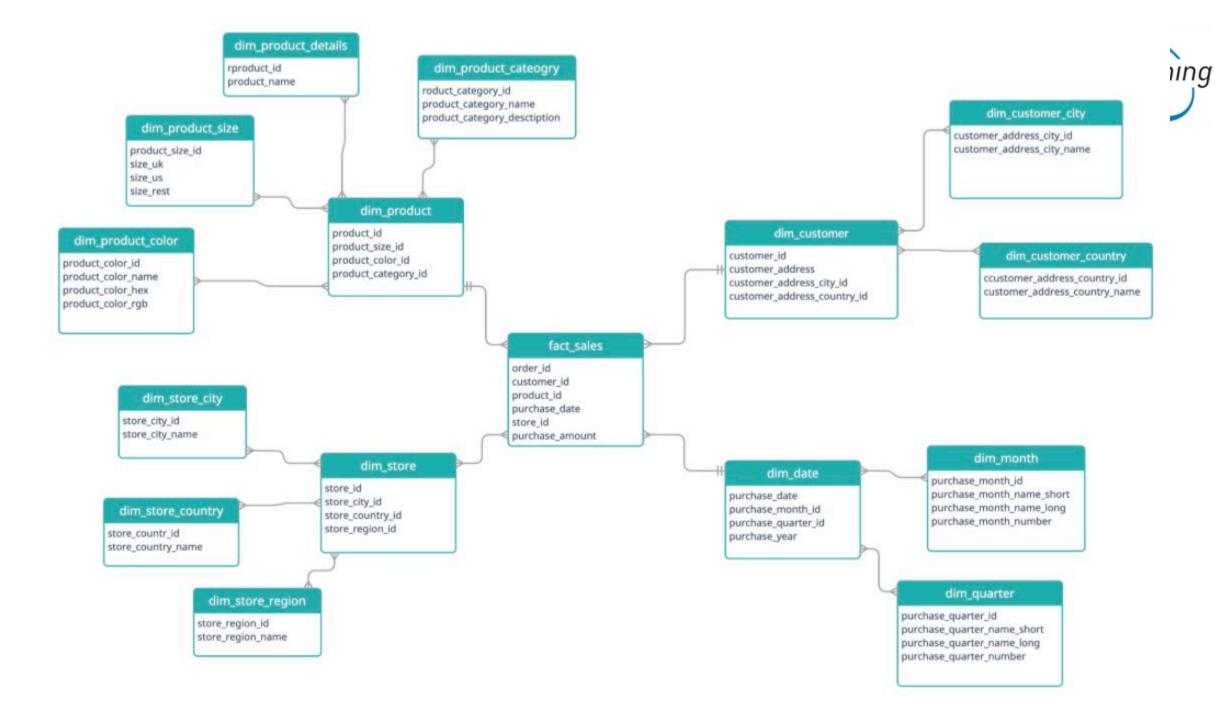




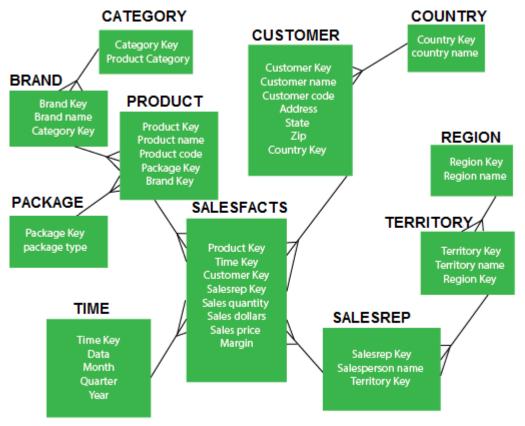






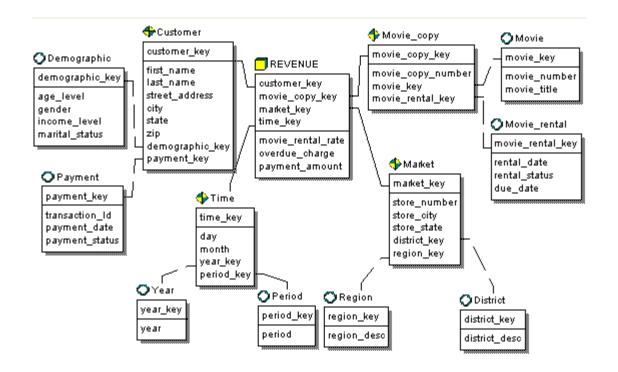


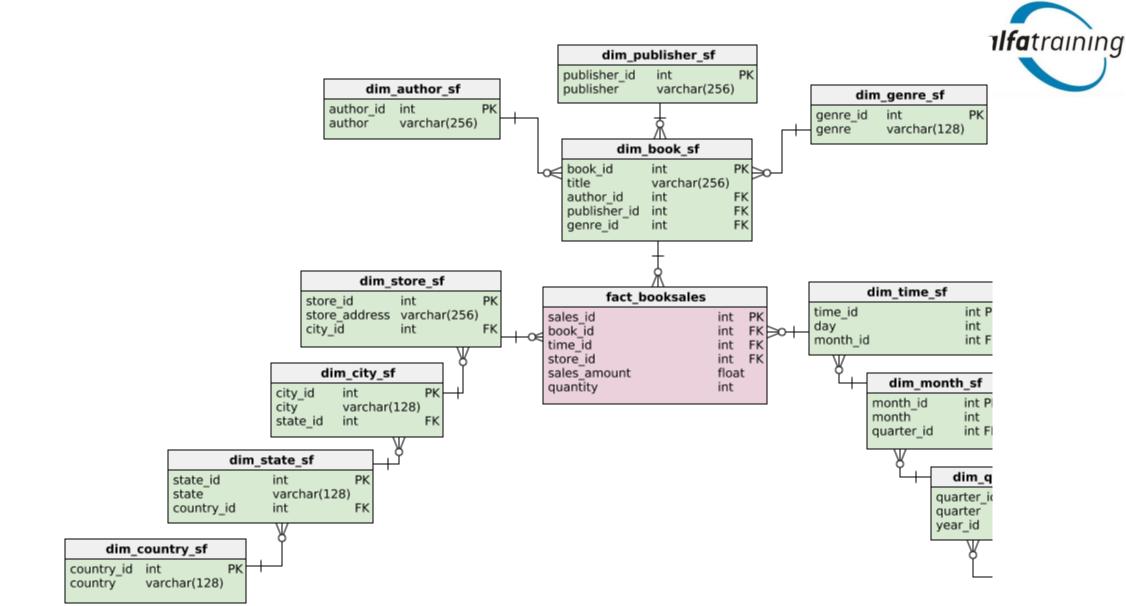




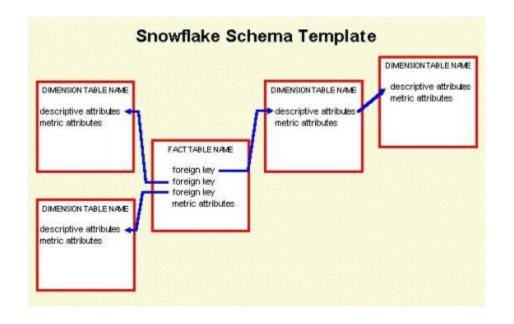
**Snowflake Schema** 





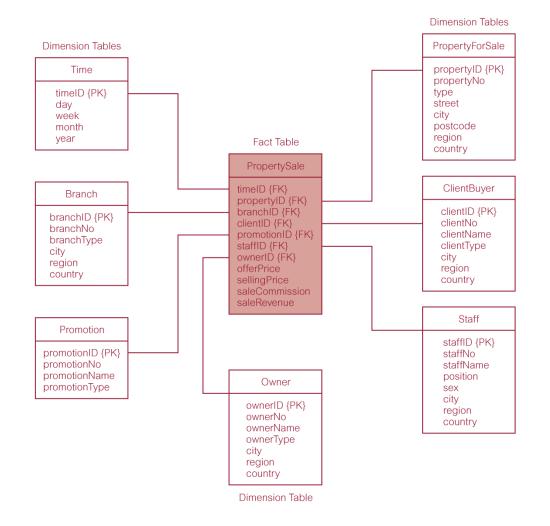






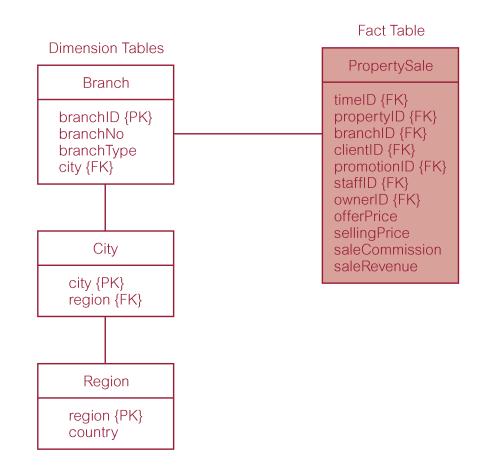
## Star schema (dimensional model) for property sales of *DreamHome*





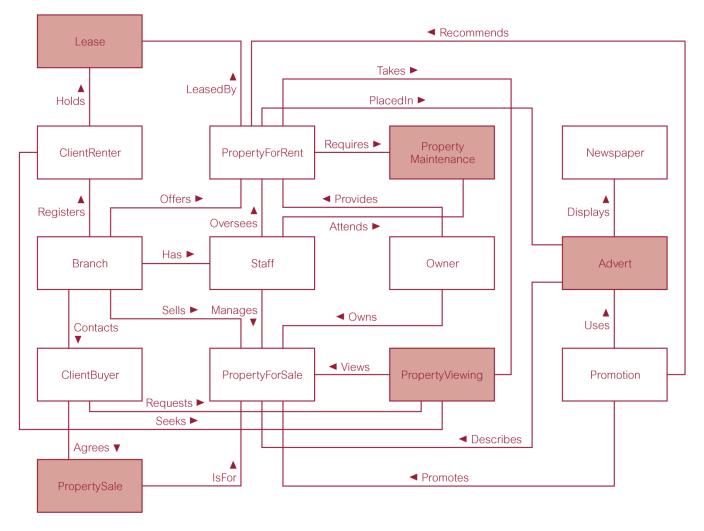


### Property sales with normalized version of Branch dimension table



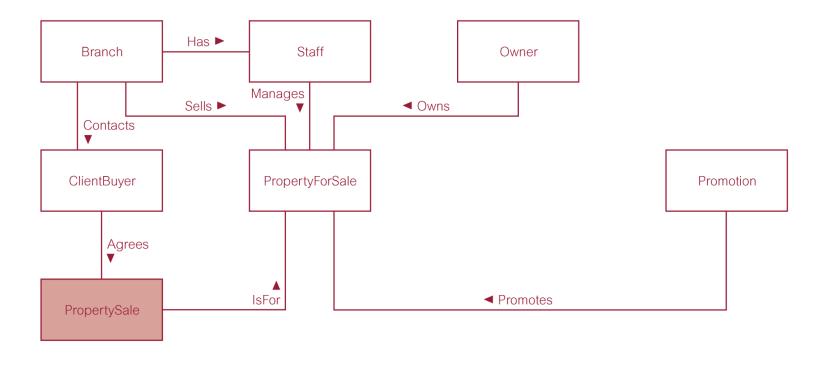
### ER model of an extended version of *DreamHome*





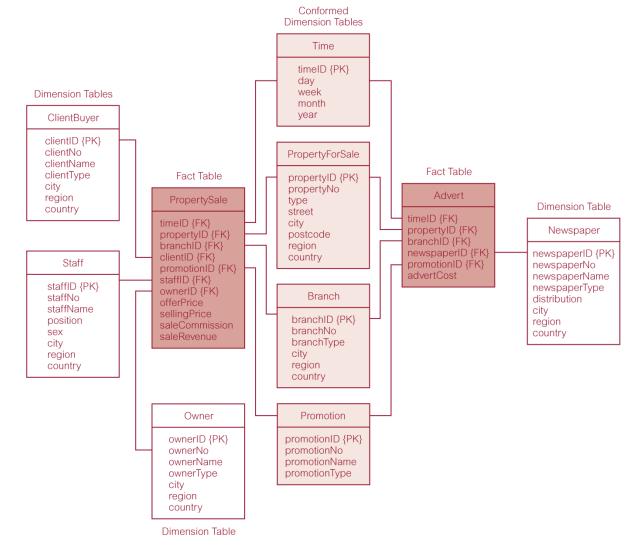


## ER model of property sales business process of *DreamHome*



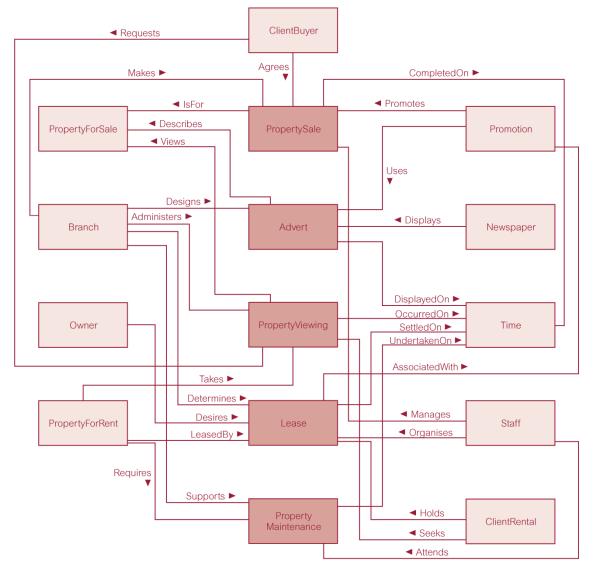
#### Star schemas for property sales and affatraining property advertising











#### Vergleich



Kriterium	Star-Schema	Snowflake-Schema
Struktur	Einfacher und denormalisierter. Ein Faktentabelle in der Mitte, umgeben von Dimensionstabellen.	Komplexer und normalisierter. Ein Faktentabelle in der Mitte, umgeben von Dimensionstabellen, die in kleinere Untertabellen aufgeteilt sind.
Performanz	In der Regel schneller, weil weniger Joins benötigt werden.	Kann langsamer sein, da mehr Joins benötigt werden, um auf die notwendigen Daten zuzugreifen.
Speicherplatz	Benötigt mehr Speicherplatz, da Daten denormalisiert sind und Duplikate existieren können.	Benötigt weniger Speicherplatz, da Daten normalisiert sind und Duplikate reduziert werden.
Einfachheit der Abfragen	Einfacher zu verstehen und zu schreiben, da weniger Joins benötigt werden.	Komplexer, da mehr Joins und eine tiefere Kenntnis der Struktur erforderlich sind.
Flexibilität	Geringer, da Änderungen an den Dimensionstabellen die Faktentabelle beeinflussen können.	Höher, da Änderungen an den Untertabellen die Haupttabellen weniger wahrscheinlich beeinflussen.
Datenaktualisierung	Einfacher, da die Daten weniger normalisiert sind.	Komplexer, da die Daten stärker normalisiert sind.

#### Vergleich



https://youtu.be/Jir31Ti8jps



# Fragen?